

NBER WORKING PAPER SERIES

THE THEORY OF PUBLIC ENFORCEMENT OF LAW

A. Mitchell Polinsky
Steven Shavell

Working Paper 11780
<http://www.nber.org/papers/w11780>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
November 2005

To appear in: A. Mitchell Polinsky & Steven Shavell (eds.), *Handbook of Law and Economics*, Volume 1, forthcoming 2006. Stanford Law School and National Bureau of Economic Research; and Harvard Law School and National Bureau of Economic Research. Research on this article was supported by the John M. Olin law and economics programs at Stanford Law School and Harvard Law School. An earlier version of portions of this chapter appeared in Polinsky and Shavell (2000a). We are at work on a book-length treatment of public enforcement of law. The views expressed herein are those of the author(s) and do not necessarily reflect the views of the National Bureau of Economic Research.

©2005 by A. Mitchell Polinsky and Steven Shavell. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

The Theory of Public Enforcement of Law
A. Mitchell Polinsky and Steven Shavell
NBER Working Paper No. 11780
November 2005
JEL No. D23, D62, D63, H23, H26, K14, K42, L51

ABSTRACT

This chapter of the forthcoming Handbook of Law and Economics surveys the theory of the public enforcement of law – the use of governmental agents (regulators, inspectors, tax auditors, police, prosecutors) to detect and to sanction violators of legal rules. The theoretical core of our analysis addresses the following basic questions: Should the form of the sanction imposed on a liable party be a fine, an imprisonment term, or a combination of the two? Should the rule of liability be strict or fault-based? If violators are caught only with a probability, how should the level of the sanction be adjusted? How much of society’s resources should be devoted to apprehending violators? We then examine a variety of extensions of the central theory, including: activity level; errors; the costs of imposing fines; general enforcement; marginal deterrence; the principal-agent relationship; settlements; self-reporting; repeat offenders; imperfect knowledge about the probability and magnitude of sanctions; corruption; incapacitation; costly observation of wealth; social norms; and the fairness of sanctions.

A Mitchell Polinsky
Stanford Law School
Stanford University
Stanford, CA 94305
and NBER
polinsky@stanford.edu

Steven Shavell
Harvard Law School
1575 Massachusetts Avenue
Hauser Hall 508
Cambridge, MA 02138
and NBER
shavell@law.harvard.edu

1. INTRODUCTION

Public enforcement of law — the use of governmental agents to detect and to sanction violators of legal rules — is a subject of obvious importance. Police and prosecutors endeavor to solve crimes and to punish criminals, regulators attempt to control violations of environmental, safety, consumer protection, and financial disclosure laws, and agents of the Internal Revenue Service seek to enforce the tax code.

The earliest economically-oriented writing on the subject of public law enforcement dates from the eighteenth century contributions of Montesquieu (1748), Beccaria (1767), and, especially, Bentham (1789). Curiously, after Bentham, the subject of law enforcement lay essentially dormant in economic scholarship until the late 1960s, when Becker (1968) published a highly influential article. Since then, several hundred articles have been written on the economics of law enforcement.¹

The main purpose of our chapter is to present the economic theory of public law enforcement in a systematic and comprehensive way.² The theoretical core of our analysis addresses the following basic questions: Should the form of the sanction imposed on a liable party be a fine, an imprisonment term, or a combination of the two? Should the rule of liability be strict or fault-based? If violators are caught only with a probability, how should the level of the sanction be adjusted? How much of society's resources should be devoted to apprehending violators?

The chapter is outlined as follows. We begin in section 2 by considering the rationale for

¹ See, for example, the references cited in Bouckaert and De Geest (1992, pp. 504-526), Garoupa (1997), Mookherjee (1997), and Polinsky and Shavell (2000a).

² For other surveys of the theory of public enforcement, see Garoupa (1997) and Mookherjee (1997). For surveys of empirical research on law enforcement and crime, see Eide (2000) and Levitt and Miles (2006), and for a survey of empirical research on enforcement of environmental regulation, see Cohen (1999, pp. 78-95).

public enforcement of law, that is, by asking why society cannot rely exclusively on private enforcement of law to control undesirable behavior. We then state the problem of public enforcement of law in general terms in section 3. In sections 4 through 6, we analyze strict liability and fault-based liability when enforcement is certain, first considering monetary sanctions, next non-monetary sanctions, and then the two together. In sections 7 through 9, we perform the same analysis when enforcement is uncertain because it is costly. In section 10 we summarize the theory, and in section 11 we discuss enforcement practices in the light of the theory we have reviewed.

We subsequently examine a variety of extensions of our core analysis in sections 12 through 28. These concern mistake, the costs of imposing sanctions, marginal deterrence, the settlement process, self-reporting of violations, corruption of law enforcers, and the fairness of sanctions, among other topics. We conclude in section 29.

2. WHY PUBLIC ENFORCEMENT RATHER THAN PRIVATE ENFORCEMENT?

Before proceeding, we should comment on the rationale for public, as opposed to private, law enforcement, where, by the latter, we mean the bringing of suits by victims of harm or those threatened by harm. An important element of the justification for private enforcement concerns information about the identity of violators. When victims of harm naturally possess knowledge of who injured them, allowing private suits for harm will motivate victims to initiate legal action and thus will harness the information they have for purposes of law enforcement. This may help to explain why, for example, the enforcement of contract law and tort law is primarily private in nature: a victim of a contract breach obviously knows who committed the breach, and a victim of a tort usually knows who the tortfeasor was. When, however, victims cannot easily identify who injured them, it may be desirable for public enforcement to be employed.

For public enforcement to be preferred in such circumstances, one still needs to explain why society cannot rely on rewards of some type to private parties other than victims (such as friends of violators or entrepreneurial private enforcers) to supply information and otherwise help in detecting violators. A difficulty with reliance on private enforcement of this sort is that if a reward is available to everyone, there might be wasteful effort devoted to finding violators (akin to excessive effort to catch fish from a common pool). Another problem is that private parties may find it hard to capture fully the benefits of developing expensive, but socially worthwhile, information systems to aid enforcement (such as computerized databases of fingerprint records). An additional obstacle to private enforcement is that force may be needed to gather information, capture violators, and prevent reprisal, yet the state frequently, if not usually, will not want to permit private parties to use force. For the preceding reasons, public enforcement often will be favored when effort is required to identify and apprehend violators.³

3. THE GENERAL PROBLEM OF PUBLIC LAW ENFORCEMENT

The general problem of public law enforcement may be viewed as one of maximizing social welfare. By social welfare, we refer to the benefits that individuals obtain from their behavior, less the costs that they incur to avoid causing harm, the harm that they do cause, the cost of catching violators, and the costs of imposing sanctions on them (including any costs associated with risk aversion). We will be explicit below about the definition of social welfare in the various contexts that we consider.

³ The differences between public and private enforcement have been discussed by Becker and Stigler (1974), Landes and Posner (1975), and Polinsky (1980a); see also Shavell (1993) and Friedman (1995). In this chapter, we assume for simplicity that public enforcement is the exclusive means of enforcement, even though in practice private parties sometimes play a complementary role by supplying information to enforcement authorities and by bringing private suits. We also abstract from private parties' efforts to protect themselves from harm (and how such efforts might relate to public enforcement), though we mention this issue in the conclusion.

The state has four major policy choices to make in undertaking law enforcement. One is about the sanctioning rule. The rule could be *strict* in the sense that a party is sanctioned whenever he has been found to have caused harm (or expected harm). Alternatively, the rule could be *fault-based*, meaning that a party who has been found to have caused harm is sanctioned only if he failed to obey some standard of behavior or regulatory requirement.

A second choice of the state concerns the form of the sanction: monetary versus non-monetary (both may be employed together). We focus on imprisonment as the primary type of non-monetary sanction and we assume that monetary sanctions are socially less costly to employ than imprisonment.

A third choice involves the magnitude of the sanction.

And the fourth choice concerns the probability of detecting offenders and imposing sanctions. This probability depends on the public resources devoted to finding violators and proving that they are liable.⁴

A. BASIC THEORY WHEN ENFORCEMENT IS CERTAIN

4. MONETARY SANCTIONS

In this section we analyze the optimal magnitude of monetary sanctions — which we call fines — assuming that enforcement is certain. We consider the two basic forms of liability, beginning with strict liability.

Suppose that individuals would obtain a gain from committing a harmful act, where the gain varies among them. If an individual does commit the act, he will have to pay a fine because

⁴ The framework for studying public law enforcement employed in this chapter derives, in many respects, from Bentham (1789). Becker (1968) first stated the enforcement problem in formal economic terms and added the choice of the probability of detection to Bentham's expression of the enforcement problem.

he is strictly liable. Let

g = gain an individual obtains if he commits the harmful act;

$z(g)$ = density of gains among individuals;

h = harm caused by an individual if he commits the harmful act;⁵

f = fine; and

w = level of wealth of an individual.

The gain could be a literal benefit, for instance, the utility obtained from taking food from a cabin in the wilderness, or the savings from not investing in some precaution, such as not obtaining pollution control equipment.⁶ The fine cannot exceed an individual's level of wealth, which is assumed to be the same for everyone.

An individual will commit the harmful act if and only if his gain from doing so exceeds the fine:⁷ $g > f$.

Social welfare equals the gains individuals obtain from committing the harmful act less the harm caused.⁸ Social welfare is not directly affected by the imposition of fines because the payment of a fine is assumed to be a socially costless transfer of money.⁹ Since the individuals

⁵ If the harm is uncertain, h can be interpreted as expected harm instead of actual harm; see section 12.

⁶ For simplicity, we assume that the gain is not itself available to pay the fine. This assumption would hold if the gain were non-monetary, as in the example of the utility benefit from taking food from a cabin. In many circumstances, the assumption would not be fitting, but the complications introduced by considering how the gain enhances an individual's ability to pay a fine would be distracting.

⁷ We assume for simplicity that he does not commit the harmful act if he is indifferent.

⁸ Some writers have questioned whether gains from committing harmful acts should necessarily be credited in social welfare; see, for example, Stigler (1970, p. 527) and Lewin and Trumbull (1990). If the gains from some type of harmful conduct were excluded from social welfare, the main consequence for our analysis would be that, for this type of conduct, society would want to achieve greater, possibly complete, deterrence. That, in turn, would tend to make a higher sanction and a higher probability of detection desirable.

⁹ In practice, of course, some costs are incurred in collecting fines, such as the cost of identifying and confiscating the individual's assets if the individual resists paying the fine. We discuss the implications of such costs in section 16.

who commit the harmful act are those whose gains exceed the fine, social welfare is

$$\int_f^{\infty} (g - h)z(g)dg. \quad (1)$$

The enforcement authority's problem is to maximize social welfare by choosing the fine f . We use an asterisk to denote the optimal value of the fine (and other variables below). It is clear from (1) that

$$f^* = h, \quad (2)$$

assuming that such a fine is feasible. Hence, individuals commit the harmful act if and only if their gain exceeds the harm, which is first-best behavior. Note, however, that there will be underdeterrence if the individual's wealth is less than h , in which case the optimal fine equals the individual's wealth.

Now consider the fault-based sanctioning rule, in which an individual who causes harm is sanctioned only if he failed to obey some standard of behavior. In the present framework, we assume that if an individual commits a harmful act, his gain must equal or exceed some threshold level of gain in order for him to escape liability; otherwise, he is considered to be at fault. Let

\hat{g} = threshold level of gain under the fault-based sanctioning rule.

If an individual's gain is less than \hat{g} he will commit the harmful act whenever $g > f$, while if his gain equals or exceeds \hat{g} , he will commit the harmful act regardless of his gain. In other words, the individual will commit the harmful act under fault-based liability if $g > \min(f, \hat{g})$.

Thus, social welfare is

$$\int_{\min(f, \hat{g})}^{\infty} (g - h)z(g)dg. \quad (3)$$

The enforcement authority's problem is to maximize (3) by choosing the standard \hat{g} and the fine f . It is straightforward to see that

$$f^* = \hat{g}^* = h \quad (4)$$

results in first-best behavior. In particular, if $g < h$, the individual would be at fault if he committed the harmful act, and would bear a fine equal to h , so he would not commit it. If $g \geq h$, he would not be at fault if he committed the harmful act and therefore would not have to pay anything, so he would commit the act. Note that f^* and \hat{g}^* are not unique: first-best behavior also can be achieved by setting $\hat{g} > h$, with $f = h$, or by keeping $\hat{g} = h$ and setting $f > h$.

Although we have seen that the first-best outcome can be achieved under both strict liability and fault-based liability (when a fine equal to h is feasible), the two forms of liability differ in the information required to implement them. To apply strict liability, the state only needs to know the harm. Under fault-based liability, however, the state needs to know more: it also needs to ascertain the gain of the individual (to determine whether he was at fault).

To illustrate strict and fault-based liability, suppose that a firm contemplates whether to discharge a pollutant that would cause harm h , rather than to transport it to a waste disposal site at cost g . Under strictly-imposed fines, the firm would incur a fine of h if it pollutes, so would not do so unless its savings from not transporting the waste, g , is greater than h . Under fault-based fines, a polluting firm would incur a fine of h only if the cost of transporting the waste g is less than h ; thus, the firm would pollute only if $g \geq h$.

5. NONMONETARY SANCTIONS

In this section we analyze the optimal magnitude of nonmonetary sanctions — which we assume to be prison sentences¹⁰ — when enforcement is certain. Let

¹⁰ It will be clear that most of what we have to say about imprisonment here and in subsequent sections would carry over, with only slight modification, to other forms of non-monetary sanctions, such as probation, electronic monitoring, community service, or, in the extreme, the death penalty.

s = length of prison sentence, where s is in $[0, s_M]$; and

$d(s)$ = disutility from prison sentence of length s ; $d(0) = 0$; $d'(s) > 0$.

The maximum sentence s_M can be interpreted as life imprisonment, and is assumed to be the same for everyone.

Under strict liability, if an individual commits the harmful act, he will bear a prison sentence.¹¹ Thus, he will commit the act if and only if his gain from doing so exceeds the disutility of the prison sentence: $g > d(s)$.

Social welfare equals the gains individuals obtain from committing the harmful act, less the harm caused, and less the cost of imposing prison sentences. The cost of prison sentences is the sum of the disutility suffered by the sanctioned individuals and the cost to the state of maintaining a prison system. Let

c = cost to the state per unit of prison sentence; $c > 0$.

Thus, social welfare is

$$\int_{d(s)}^{\infty} [g - h - (d(s) + cs)]z(g)dg. \quad (5)$$

The enforcement authority's problem is to maximize social welfare by choosing the prison sentence s . From (5), the first-order condition can be written as

$$(h + cs)(dZ(d(s))/ds) = \int_{d(s)}^{\infty} [d'(s) + c]z(g)dg, \quad (6)$$

where $Z(\cdot)$ is the cumulative distribution function for $z(\cdot)$. The left-hand side is the marginal benefit of raising s : $h + cs$ is the social gain from deterring the marginal individuals, while

¹¹ Strict liability is in fact an unusual form of liability when the sanction is imprisonment. (An example of a strict liability criminal offense is the serving of liquor by a bar to underage individuals.) As will be seen, fault-based liability tends to be superior to strict liability when the sanction is imprisonment.

$dZ(d(s))/ds$ is the number of such individuals deterred.¹² The right-hand side is the marginal cost of raising s : individuals who are not deterred incur additional disutility $d'(s)$ and cause the state to incur additional imprisonment costs c .

The optimal prison sentence could lead to either underdeterrence or overdeterrence due to the costs of prison sentences. To see why, suppose the sentence were such that individuals committed the harmful act if and only if their gain equaled the harm. If the sentence were raised above this level, some individuals would be deterred, which would reduce prison costs, but those who would not be deterred would bear longer sentences. Depending on which effect is stronger, it may be desirable to raise, or lower, the sentence, leading to overdeterrence or underdeterrence. Note that a marginal change of the sentence from the initial level does not affect the gains net of the harm because the marginal individuals are those whose gains equal the harm.

Regardless of whether the optimal prison sentence causes underdeterrence or overdeterrence, the strict sanctioning rule does not achieve the first-best outcome because it leads to the imposition of costly sanctions.

Now consider the fault-based sanctioning rule, with a standard \hat{g} . Analogously to the case of fines, an individual will commit the harmful act under fault-based liability if $g > \min(s, \hat{g})$. If $s < \hat{g}$, then there will be some individuals who commit the harmful act who will be found at fault, those for whom $s < g < \hat{g}$. Then social welfare would be

$$\int_s^{\infty} (g - h)z(g)dg - \int_s^{\hat{g}} (d(s) + cs)z(g)dg. \quad (7)$$

If $s \geq \hat{g}$, then social welfare is

¹² The disutility of imprisonment and the benefit from committing the harmful act do not appear in (6) because they offset each other for the marginal person.

$$\int_{\hat{g}}^{\infty} (g - h)z(g)dg. \quad (8)$$

The enforcement authority's problem is to maximize (7) or (8) by choosing the standard \hat{g} and the sentence s . It is straightforward to see that

$$s^* = \hat{g}^* = h \quad (9)$$

results in the first-best outcome. First observe that first-best behavior will occur. In particular, if $g < h$, the individual would be at fault if he committed the harmful act, and bear a sentence equal to h , so he would not commit it. If $g \geq h$, he would not be at fault if he committed the harmful act and therefore would not bear any sentence, so he would commit the act. Second, given (9), costly sanctions will not be imposed because individuals will never choose to be at fault. Hence, the first-best outcome is achieved. Note that s^* and \hat{g}^* are not unique: the first-best outcome also can be achieved by keeping $\hat{g} = h$ and setting $s > h$.

The preceding discussion shows that when sanctions are costly, the fault-based sanctioning rule is superior to the strict sanctioning rule. Not only does the fault-based rule lead to first-best deterrence, it does so without anyone actually incurring a costly sanction.¹³

6. COMBINED SANCTIONS

In this section we consider the optimal mix of fines and imprisonment sanctions when they can be used together. Under the strict sanctioning rule, social welfare in this case is

$$\int_{f+d(s)}^{\infty} [g - h - (d(s) + cs)]z(g)dg. \quad (10)$$

First observe that it cannot be optimal to employ a prison sentence unless the fine has been set as high as possible, equal to individuals' wealth level w . To see why, suppose

otherwise, that $f < w$ and $s > 0$. Then it would be possible to raise f and lower s so as to keep $f + d(s)$ constant, thereby leaving behavior unaffected but lowering the cost of imprisonment (see (10)).

Whether it is optimal to use a prison sentence in addition to a fine depends on whether a fine alone is sufficient to achieve first-best deterrence. If wealth levels are high enough, if $w \geq h$, then a fine equal to harm is feasible and there would be no need for a prison sentence. If, however, $w < h$, then relying on fines alone would lead to underdeterrence and it might be desirable to employ a prison sentence in addition to the maximal fine of w . To see whether a prison sentence would be desirable in these circumstances, consider the derivative of (10) with respect to the sentence s when $f = w$:

$$[h - w + cs][dZ(w + d(s))/ds] - \int_{w + d(s)}^{\infty} (d'(s) + c)z(g)dg. \quad (11)$$

It follows that the condition for a positive prison sentence to be optimal is that

$$[h - w][dZ(w)/ds] > \int_w^{\infty} (d'(0) + c)z(g)dg. \quad (12)$$

The left-hand side of (12) is the value of deterring the marginal individuals, while the right-hand side is the marginal cost of imposing prison sentences on individuals who are not deterred. If a positive prison sentence is socially desirable, it is determined from the first-order condition associated with (11).

Next consider the optimal mix of fines and imprisonment under the fault-based sanctioning rule. The key point is that, unlike under strict liability, it is always desirable to employ prison sentences to obtain compliance with the fault standard if fines alone are inadequate to do so. This is because, as noted above, sanctions are not actually imposed when

¹³ But see section 15 below on the possibility of error and thus of the bearing of sanctions.

the fault standard is complied with, so there is no social cost from using the threat of prison sentences to obtain compliance. Thus, any combination of fines and prison sentences that induces potential offenders to comply with the fault standard is optimal.

B. BASIC THEORY WHEN ENFORCEMENT IS UNCERTAIN

In this part we investigate the level of enforcement resources that the state should devote to detecting offenders. We assume that the higher the level of expenditures on enforcement, the greater is the probability of detection. Let

e = enforcement expenditures by the state; and

$p(e)$ = probability of detection given e ; $p'(e) > 0$; $p''(e) < 0$.

We derive the optimal probability of detection¹⁴ and, along with it, the optimal magnitude of sanctions.¹⁵

7. MONETARY SANCTIONS

In this section we analyze the optimal probability and magnitude of fines, first assuming that individuals are risk neutral and then that they are risk averse.

7.1 The Risk-Neutral Case

Under the strict liability rule, an individual will commit the harmful act if and only if his

¹⁴ We implicitly assume that enforcement expenditures e determine a single probability of detection. Thus, we do not consider an issue identified by Lando and Shavell (2004), that it may be advantageous to concentrate enforcement resources on a subset of potential offenders (for example, by auditing taxpayers whose last names begin with certain letters) rather than to spread enforcement resources evenly.

¹⁵ Although we assume that the probability of detection can be set independently of the level of sanctions, the two might be connected. This is because high sanctions may lead juries to be less likely to convict defendants, or may induce individuals to engage in greater efforts to avoid detection; on these points, see Andreoni (1991) and Malik (1990), respectively. See also Bar-Gill and Harel (2001) for a discussion of how the level of crime affects both the probability and magnitude of sanctions.

gain from doing so exceeds the expected fine: $g > pf$. Social welfare, which now reflects the enforcement expenditures of the state, is

$$\int_{p(e)f}^{\infty} (g - h)z(g)dg - e. \quad (13)$$

The enforcement authority's problem is to maximize (13) by choosing enforcement expenditures e (and thus the probability of detection p), as well as the level of the fine f .

Before considering the complete problem, suppose that enforcement expenditures are fixed, resulting in the probability of detection p . It is obvious that if $pf = h$, namely, the expected fine equals the harm, (13) will be maximized over f since the first-best outcome will be achieved. In other words,

$$f^* = h/p, \quad (14)$$

assuming that such a fine is feasible. Individuals then commit the harmful act if and only if their gain exceeds the harm, which is first-best behavior.¹⁶ Note, however, that there will be underdeterrence if individuals' wealth is less than h/p , in which case the optimal fine equals their wealth.

Now suppose that both enforcement expenditures and the fine are chosen by the state.¹⁷ Then the optimal fine is maximal: $f^* = w$. To demonstrate this, suppose that f is less than w . Then f can be raised and e lowered so as to keep $p(e)f$ — the level of deterrence — constant. Because the behavior of individuals is unaffected but enforcement expenditures fall, social welfare rises (the first term in (13) does not change but e is lower). Hence, the optimal f cannot be less than w . In other words, because any particular level of deterrence can be achieved with

¹⁶ The general formula (14), or its equivalent, was put forward by Bentham (1789, p. 173), was emphasized by Becker (1968), and has been noted by many others since then.

¹⁷ Consideration of this issue originated with Becker (1968); as we noted, early writers on law enforcement did not examine the issue of the choice of enforcement effort.

different combinations of the fine and the probability of detection, society should employ the highest possible fine and a correspondingly low probability of detection in order to economize on enforcement expenditures.¹⁸

We next show that the optimal probability of detection is such that the expected fine is less than the harm, $p(e^*)w < h$ — that is, some degree of underdeterrence is desirable. Observe that the first-order condition determining optimal enforcement expenditures e^* is

$$[h - p(e)w][dZ(p(e)w)/de] = 1. \quad (15)$$

The left-hand side is the marginal social benefit of the deterrent effect from a higher probability of detection. The right-hand side is the marginal cost of greater spending on enforcement. It follows from (15) that $p(e^*)w < h$. To explain this result, suppose that p were such that $pw = h$. Then there would be no first-order loss of social welfare from lowering p because the individuals who would be induced to engage in the harmful activity would obtain gains equal to harm. But enforcement costs would be saved, making it desirable to lower the probability. How much p should be lowered depends on the resulting savings in enforcement expenses compared to the net social costs of underdeterrence.¹⁹

Under fault-based liability, analogues of the above conclusions hold. The optimal fine is maximal and the optimal probability of detection is such that the expected fine is less than the harm. Moreover, the optimal fault standard is less than the first-best standard. The explanation

¹⁸ Although the general point that a low probability-high fine combination conserves enforcement costs is due to Becker (1968), he did not formally consider bounds on fines (and much of his analysis implicitly presumes that the optimal fine is not maximal). Carr-Hill and Stern (1979, pp. 300-304) and Polinsky and Shavell (1979, pp. 883-884) observed that Becker's argument implies that the optimal fine is equal to its upper bound. Many scholars have noted the unrealism of this result and have introduced additional considerations that imply that less-than-maximal fines are optimal. We will discuss several important factors of this type, including risk aversion, general enforcement, and marginal deterrence. See also Andreoni (1991), Bebchuk and Kaplow (1992; 1993), Malik (1990), and Polinsky and Shavell (1991;200b) for discussion of other such considerations.

¹⁹ The point of this paragraph — that some underdeterrence is optimal — was first made by Polinsky and Shavell (1984).

for these results is essentially that given above. The fine is maximal in order to reduce the probability of detection and thereby save enforcement costs, and some underdeterrence is desirable because this also allows savings in enforcement costs.

7.2 The Risk-Averse Case

Now suppose that individuals are risk averse and that social welfare is the sum of expected utilities of individuals (in the risk-neutral case, social welfare was equivalent to the sum of utilities). For convenience, assume that the risk of being harmed is the same for everyone and that individuals buy insurance against harm, paying a premium equal to the expected harm. An individual's wealth is his initial wealth, less the taxes he pays, less the expected harm he suffers, and less the fine if he has to pay it. Let

$U(\cdot)$ = utility of wealth; U is concave in wealth;

t = tax; and

λ = fraction of population that commits the harmful act.

Both t and λ are endogenous.

An individual will commit the harmful act if $g + (1 - p)U(w - t - \lambda h) + pU(w - t - \lambda h - f)$ is greater than $U(w - t - \lambda h)$, or equivalently, if

$$g > p[U(w - t - \lambda h) - U(w - t - \lambda h - f)]. \quad (16)$$

Note that we are implicitly assuming that the gain g is non-monetary.²⁰ The condition (16) implicitly determines the fraction of the population λ that commits the harmful act. The tax t is such that the government breaks even; hence t equals the enforcement cost e less the fine revenue collected λpf .

Social welfare, the sum of individuals' expected utilities, equals

²⁰ If the gain were monetary, then the condition would become $(1 - p)U(w + g - t - \lambda h) + pU(w + g - t - \lambda h - f) > U(w - t - \lambda h)$. The qualitative nature of the results in this section would not be affected.

$$(1 - \lambda)U(w - t - \lambda h) + \int_{p[U(w - t - \lambda h) - U(w - t - \lambda h - f)]}^{\infty} [g + (1 - p)U(w - t - \lambda h) + pU(w - t - \lambda h - f)]z(g)dg - e, \quad (17)$$

since the individuals who commit the harmful act are those whose gains exceed the expected disutility of the fine.

Suppose initially that the probability of detection is fixed. Then the optimal fine in the risk-averse case tends to be lower than that in the risk-neutral case for two reasons. First, lowering the fine reduces the bearing of risk by individuals who commit the harmful act. Second, because risk-averse individuals are more easily deterred than risk-neutral individuals, the fine does not need to be as high to achieve any desired degree of deterrence.²¹

Now consider choosing both the probability and magnitude of fines. The optimal fine generally is not at its maximum when individuals are risk averse. This is because the use of a very high fine would impose a substantial risk-bearing cost on individuals who commit the harmful act. More precisely, reconsider the argument employed in the risk-neutral case. If f is less than the maximal fine (now $w - t - \lambda h$), it still is true that f can be raised and e lowered so as to keep deterrence constant. But due to risk aversion, the probability of detection that maintains deterrence falls *more* than proportionally, implying that the expected fine, and therefore fine revenue, falls. This reduction in fine revenue reflects the disutility caused by imposing greater risk on risk-averse individuals. If individuals are sufficiently risk averse, the decline in fine revenue associated with greater risk bearing could more than offset the savings in enforcement expenditures from reducing the probability of detection, implying that social welfare would be

²¹ It is possible, however, that the optimal fine is higher in the risk-averse case than in the risk-neutral case, for the following reason. A way to reduce the bearing of risk is to deter more individuals from committing the harmful act, for then fewer individuals will be subject to the risk of the fine. See Kaplow (1992).

lower.²²

In effect, when individuals are risk averse, fines become a socially costly sanction (reflected in an increase in taxes) rather than a mere transfer of wealth. The more risk averse individuals are, the better it is to control their behavior by using a lower fine and a higher probability of detection, even though this raises enforcement costs.

As in the risk-neutral case, there is a reason when individuals are risk averse to reduce enforcement costs by setting the probability such that some individuals will commit the harmful act even though their gain is less than the harm —meaning that there will be some underdeterrence.²³

Under fault-based liability, the conclusions are different from those under strict liability. The differences are due to the fact that individuals are induced to comply with the fault standard and therefore do not bear risk. Consequently, the results are the same as those in the risk-neutral case: the optimal fine is maximal; the optimal probability is relatively low; and the optimal fault standard is such that there is some underdeterrence.

Because, as just emphasized, fines are not actually imposed under fault-based liability, fault-based liability is superior to strict liability, under which risk is borne. This advantage of fault-based liability is analogous to the advantage of fault-based liability over strict liability when imprisonment is used — a costly sanction is not actually imposed.

²² The point that the optimal fine may be less than maximal when individuals are risk averse was proved initially by Polinsky and Shavell (1979) in a model with two levels of gain. See also Kaplow (1992), who demonstrates in an example that the fine may be less than maximal. It can be shown in the general model under discussion here that the optimal fine must be less than maximal if the cost of raising the probability of detection is sufficiently small (given that the wealth of individuals is not too low). The idea of the proof is that, if the cost of raising the probability were zero, the optimal probability would be one and the optimal fine less than maximal (equal to the harm). The conclusion follows by a continuity argument.

²³ It also is possible, however, that overdeterrence would be optimal. The reason is that the imposition of risk can be reduced by discouraging individuals from engaging in the harmful activity.

Another advantage of fault-based liability is that it may result in lower enforcement expenditures than under strict liability. Specifically, because fines are not imposed under fault-based liability, it becomes desirable to use a high (maximal) fine, which allows a relatively low probability of detection to be employed.

Of course, in reality, as we will be discussing below, mistakes may occur under fault-based liability, resulting in the imposition of risk. To the extent that risk exists under fault-based liability, the main result obtained under strict liability — that fines should not be maximal — carries over to fault-based liability.

8. NONMONETARY SANCTIONS

In this section we analyze the optimal probability and magnitude of prison sentences, first assuming that individuals are risk neutral, then that they are risk averse, and finally that they are risk preferring. The last case is of particular relevance for imprisonment sanctions, as will be explained.

8.1 The Risk-Neutral Case

We assume here that $d(s) = s$, that is, that the disutility of imprisonment rises in proportion to the length of the sentence. This implies that individuals are indifferent between a sure sentence of s and an uncertain sentence with a mean of s . Thus, individuals display a risk-neutral attitude towards imprisonment sentences.

Under the strict liability rule, an individual will commit the harmful act if and only if his gain exceeds the expected sentence: $g > ps$. Social welfare is

$$\int_{p(e)s}^{\infty} (g - h - p(e)(s + cs))z(g)dg - e. \quad (18)$$

The enforcement authority's problem is to maximize (18) by choosing enforcement expenditures

e and the length of the sentence s .

In this case the optimal sentence is maximal: $s^* = s_M$. As seen above, when the sanction is a fine, the optimal sanction is maximal. This is also true here even though the sanction is costly to impose. To demonstrate that $s^* = s_M$, suppose that s is less than s_M . Then s can be raised and e lowered so as to keep $p(e)s$ constant. The behavior of individuals is unaffected because $p(e)s$ has not changed. The social cost of imprisonment also is unaffected because $p(e)(s + cs)$ is constant. In other words, although the sentence is higher, proportionally fewer individuals are imprisoned. But because enforcement expenditures fall, social welfare rises. Hence, the optimal s equals s_M .

The optimal probability of detection is such that the expected prison sentence could lead to either underdeterrence or overdeterrence. This is essentially for the reasons discussed above (see section 5) when the probability of detection was fixed at one. Here, however, there is an additional factor favoring underdeterrence, namely that by lowering the probability, enforcement resources are saved.

Under fault-based liability, first observe that it must be optimal to have compliance with the fault standard. If the expected sentence were less than the standard, so that some individuals would choose to violate the standard and bear the expected sentence, then it would be optimal to lower the standard to the expected sentence. For then, the behavior of individuals would be the same, but the cost of imprisonment would be avoided.

Next observe that the expected sentence must equal the standard, rather than be higher. Otherwise, the probability of detection could be lowered without affecting behavior, but enforcement costs would be saved.

It is clear, too, that the optimal sentence must be maximal. This is for the now familiar reason that, otherwise, the sentence could be raised and the probability of detection lowered

proportionally, without affecting behavior, but saving enforcement expenditures (imprisonment costs are zero because of compliance with the standard).²⁴ Given $s^* = s_M$, the optimal probability satisfies $p(e)_{s_M} = \hat{g}$ because of our observation in the previous paragraph.

The optimal standard \hat{g}^* is less than the harm h ; in other words, there is some underdeterrence. This is true for a reason already discussed: by lowering the standard, enforcement costs are saved, but there is no first-order effect on social welfare due to more individuals committing the harmful act because their gain equals the harm.

Finally, as we previously emphasized, fault-based liability possesses the advantage over strict liability that costly sanctions are not actually imposed (in the absence of mistakes). Moreover, because the optimal sentence is maximal for this reason, a low probability of detection can be used.

8.2 The Risk-Averse Case

We assume now that the disutility of the sentence $d(s)$ rises more than in proportion to the sentence. This could occur because of an increasing desire of a prisoner to join the outside world or a growing distaste for the prison environment as time in jail increases. This assumption implies that individuals prefer a sure sentence of s to an uncertain sentence with a mean of s . In other words, individuals are risk-averse in imprisonment sentences.

Under the strict liability rule, an individual will commit the harmful act if and only if his gain exceeds the expected disutility of the sentence, $g > pd(s)$, and social welfare is

$$\int_{p(e)s}^{\infty} (g - h - p(e)(d(s) + cs))z(g)dg - e. \quad (19)$$

The optimal sentence is again maximal: $s^* = s_M$. The reasons given in the risk-neutral

²⁴ Note, too, that the expected sentence remains constant, so even if imprisonment costs were borne, they would not change.

case are reinforced here. As s is raised and e lowered so as to keep $p(e)d(s)$ and the behavior of individuals constant, p declines proportionally more than s rises because of individuals' risk aversion. In other words, ps declines, which implies that the public cost of imprisonment pcs falls.²⁵ Hence, social welfare rises both because the cost of enforcement declines and the cost of imprisonment declines.

The optimal probability of detection is such that, as in the risk-neutral case, the expected prison sentence could lead to underdeterrence or overdeterrence.

Under fault-based liability, the results are essentially the same as in the risk-neutral case because, if the fault standard is complied with, risk is not borne. Thus, it is optimal to have compliance with the fault standard; the expected disutility of the sentence must equal the standard; the optimal sentence is maximal; there is some underdeterrence; and fault-based liability is superior to strict liability. A difference, however, is that the probability of detection needed to obtain compliance with the fault standard can be lowered because individuals are risk averse.

8.3 The Risk-Preferring Case

Finally, suppose that the disutility of the sentence $d(s)$ rises less than in proportion to the sentence. This could occur because the disutility from the stigma of being in jail might be substantial from having spent even a short amount of time there, but not increase much with the length of imprisonment.²⁶ Individuals' discounting of the future disutility of imprisonment also makes earlier years of imprisonment more important than later ones. The present assumption implies that individuals prefer an uncertain sentence with a mean of s to a sure sentence of s ;

²⁵ Note the contrast with the case of risk aversion in fines, in which the decline in the expected sanction pf meant a decline in revenue to the state and thus an increase in taxes. Here the decline in the expected sanction ps means a decline in expenses to the state and thus a decrease in taxes.

individuals are risk-preferring in imprisonment sentences.

Under the strict liability rule, social welfare is again given by (19). The optimal sentence, however, might be less than maximal: $s^* \leq s_M$. Now, when the sentence is raised, the probability that maintains deterrence cannot be lowered proportionally, implying that the expected prison term rises. Because the resulting increased cost to the public of imposing imprisonment sanctions might exceed the savings in enforcement expenditures from lowering the probability, the optimal prison term might not be maximal.²⁷

Under fault-based liability, the results are again essentially the same as in the risk-neutral case.

9. COMBINED SANCTIONS

Under the strict sanctioning rule, as we explained above, it never is optimal to employ a prison sentence unless the fine has been set as high as possible, since fines are socially cheaper sanctions. Whether it is optimal to use a prison sentence in addition to the maximal fine depends on the extent of underdeterrence that would result if fines were used alone, and the social cost of imprisonment.

Under the fault-based sanctioning rule, the key point is that it is always desirable to employ the maximal prison sentence in addition to the maximal fine, since neither sanction is actually imposed. By using maximal sanctions, the probability of detection can be set at a low level, thereby saving enforcement costs.

²⁶ Also, the first years of imprisonment may create special disutility due to brutalization of the prisoner.

²⁷ The results in this section were first presented by Polinsky and Shavell (1999) (although Shavell (1991b) notes the result in the case of risk neutrality).

C. BASIC THEORY SUMMARIZED AND COMPARED TO PRACTICE

10. SUMMARY OF THE BASIC THEORY

In this section we summarize the main points from parts A and B:

(a) When the probability of detection of a harmful act is taken as fixed and individuals are risk neutral, the optimal fine is the harm divided by the probability of detection, for this results in an expected fine equal to the harm. However, the risk aversion of individuals tends to lower the level of the optimal fine.

(b) When the probability of detection can be varied, relatively high sanctions may be optimal, for this allows a relatively low probability to be employed and thereby saves enforcement costs. Indeed, the optimal fine is maximal if individuals are risk neutral in wealth, and the optimal imprisonment term is maximal if individuals are risk neutral or risk averse in imprisonment. Optimal sanctions might not be maximal, however, when individuals are risk averse in wealth or risk preferring in imprisonment, both plausible assumptions, although the motive to set sanctions at relatively high levels in order to reduce enforcement costs still applies.²⁸

(c) Optimal enforcement tends to be characterized by some degree of underdeterrence relative to first-best behavior, because allowing underdeterrence conserves enforcement resources. More precisely, by lowering the probability of detection slightly from a level that would lead to first-best behavior, the state reduces enforcement costs, and although more individuals commit the harmful act, these individuals do not cause social welfare to decline substantially because their gains are approximately equal to the harm.

²⁸ There are other reasons why optimal sanctions might not be maximal. See, for example, the discussion of general enforcement in section 17.

(d) The use of fines should be exhausted before resort is made to the costlier sanction of imprisonment.

(e) An advantage of fault-based liability over strict liability is that sanctions that are costly to impose — imprisonment, and fines when individuals are risk averse — are imposed less often under the former rule. Under fault-based liability, individuals generally are induced (in the absence of mistakes) to obey fault standards, and therefore ordinarily do not bear sanctions. Under strict liability, however, individuals who cause harm are sanctioned whenever they are caught.

(f) An advantage of strict liability over fault-based liability is that the former is easier to apply. Strict liability requires the state to determine only the harm done, whereas fault-based liability requires the state to ascertain optimal behavior (in order to set the fault standard) and to observe actual behavior (in order to apply the standard).

11. THEORY VERSUS PRACTICE

Having reviewed the basic theory of public enforcement of law, we briefly comment on the relationship between optimal enforcement and enforcement in practice.

First observe that important features of actual public law enforcement are congruent, at least in a broad sense, with what is theoretically desirable. Public enforcement is often characterized by low probabilities of detection. This is true for many criminal acts, and also is frequently the case in other spheres of public enforcement, such as traffic control and tax collection.²⁹ That probabilities of detection are low undoubtedly reflects the cost of raising the

²⁹ U.S. Department of Justice (1997b, p. 205, table 25) indicates, for example, that the likelihood of arrest was 13.8% for burglary, 14.0% for automobile theft, and 16.5% for arson. Kenkel (1993, p. 145) estimates that the probability of arrest for drunk driving is “only about .003.” And according to Andreoni et al. (1998, p. 820), the audit rate for individual tax returns was 1.7 percent in 1995.

probability, a central factor in our discussion.

Corresponding to the low probabilities of detection are relatively high sanctions, often exceeding harm. For example, it seems that the sentence for theft typically outweighs the harm from that act, that the penalty for double parking frequently surpasses the resulting congestion costs, and that the sanction for tax evasion tends to exceed the social losses thereby created. Sanctions that are in excess of harm are needed for proper deterrence when the probabilities of enforcement are less than one, as they are in these examples.

The theory of optimal public enforcement of law also helps to explain why society uses the sanction of imprisonment when it does — for the category of harmful acts labeled criminal, notably, for theft, robbery, rape, murder, and so forth.³⁰ Because such acts cause substantial harm, yet often are detected with a low probability, the magnitudes of desirable penalties are high. If these penalties were solely monetary, they often would exceed the assets of the individuals who commit the acts, for individuals who commit crimes tend to have very low assets.³¹ Imprisonment sanctions, therefore, usually will be required to maintain an adequate level of deterrence of acts classified as criminal.³²

Note, too, that the standard of liability when imprisonment sanctions are imposed is generally fault-based — imprisonment is premised on the nature of the wrongful act, not merely on the fact that harm was done. This is socially desirable because, as we stressed, fault-based

³⁰ See generally Posner (1985, pp. 1201-1205), Shavell (1985, pp. 1236-1241), and Shavell (2004, pp. 543-568).

³¹ For example, in U.S. Department of Justice (1988, p. 35) it is reported that “the average inmate was at the poverty level before entering jail” and in U.S. Department of Justice (1998, p. 4) it is stated that almost half of jail inmates reported incomes of less than \$600 a month in the month before their most recent arrest.

³² The use of imprisonment sanctions also makes sense in view of their incapacitative function: crimes cause substantial harm and may be difficult or expensive to deter (for the reason we just emphasized, as well as others), so that it often will be desirable to incapacitate individuals who have committed them. See section 25.

liability reduces the use of socially costly sanctions.

Although actual public enforcement is consistent in many respects with the theory of optimal enforcement, actual enforcement also appears to deviate in various ways from what is theoretically desirable. We note two discrepancies of general importance. First, substantial enforcement costs could be saved without sacrificing deterrence by reducing enforcement effort and simultaneously raising fines. This is possible in many enforcement contexts because fines are presently very low relative to the assets of violators. For example, fines for most parking violations are less than \$50, penalties for underpayment of income taxes are typically on the order of 20% of the amount not paid, and fines for corporate violations of health and safety regulations are frequently minuscule in relation to corporate assets. In such areas of enforcement, therefore, fines could readily be, say, doubled and enforcement costs reduced significantly, while maintaining deterrence at present levels.

Not only can present levels of deterrence be achieved more cheaply, it also seems that these levels are often too low. This is a reasonable supposition given the limited use of fines that we just noted and the low probabilities of their application. For example, the probability of a tax audit is approximately 2%; when combined with the modest penalties for underpayment, one would predict substantial tax avoidance.³³ Evidence also suggests that the expected fine for driving while intoxicated is on the order of one-quarter of the expected harm caused by such behavior,³⁴ and that monetary sanctions imposed on corporations equal on average only thirty-

³³ In 1995 the audit rate for individual returns was 1.7 percent, as noted above, and the civil penalty for underpayment of taxes ordinarily is calculated as 20 percent of the underpayment that results from wrongful conduct (such as substantially misstating a valuation). See Andreoni et al. (1998, p. 820). Thus, for every dollar of underpayment, the expected payment, including the underpayment and the civil penalty, is only approximately \$0.02 (= .017 x \$1.20).

³⁴ See Kenkel (1993, p. 145). The expected fine is \$12.82 and the expected harm is \$47.77 (both in 1986 dollars). While the latter number may seem low, keep in mind that it is the product of the probability that a harm will occur as a result of drunken driving, and the level of harm if harm does occur. (To properly determine whether

three percent of the harms caused.³⁵ Given the ample opportunities that exist for augmenting penalties, as well as the possible desirability of increasing enforcement effort, society probably should raise levels of deterrence in many areas of enforcement.

D. EXTENSIONS OF THE BASIC THEORY

This concludes the presentation of the basic theory of public enforcement of law. We now turn to various extensions and refinements of the basic analysis.

12. ACCIDENTAL HARMS

In our analysis above, we implicitly assumed that the acts that individuals commit result in harm with certainty. In many circumstances, of course, acts result in harm only with a probability. A driver that speeds only creates a likelihood of a collision; or a firm that stores toxic chemicals in a substandard tank only creates a probability of a harmful spill.

Essentially all of the results in the basic analysis carry over in a straightforward way when harms are accidental. If individuals are risk neutral, sanctions are monetary, and the expected sanction equals harm whenever harm turns out to occur, then induced behavior will still be socially optimal; further, the optimal magnitude of sanctions is maximal if individuals are risk neutral because this allows enforcement costs to be saved, but is not necessarily maximal if individuals are risk averse, and so forth. Our general conclusions in the basic analysis can thus be interpreted to apply both when harms occur for sure and when harms occur accidentally.

dangerous driving is underdeterred, one also would have to take into account the threat of liability from private suits brought by accident victims. But the deterrent effect of such suits will be dulled to the extent that drivers do not have sufficient assets to pay for the harms suffered by accident victims, or have liability insurance and therefore only partially bear the financial consequences of a lawsuit.)

³⁵ See Cohen (1989, pp. 617-618, 658). Cohen notes, however, that he did not take into account other sanctions imposed on corporate criminals, including restitution, civil penalties, and private tort suits.

There is, however, a new issue that arises when harm is uncertain: a sanction can be imposed either on the basis of the commission of a dangerous *act* that increases the chance of harm — storing chemicals in a substandard tank — or on the basis of the *actual occurrence of harm* — only if the tank ruptures and results in a spill. In principle, either approach can be employed to achieve optimal deterrence. To illustrate, suppose that the substandard tank has a 10% chance of rupturing, in which case the harm would be \$10 million; the expected harm from using the tank therefore is \$1 million. If individuals are risk neutral and sanctions are imposed only when harm occurs, deterrence will be optimal if, as usual, the sanction equals the harm of \$10 million. Alternatively, if sanctions are imposed on the basis of the dangerous act of using the substandard tank, deterrence will be optimal if the owner of the tank faces a sanction equal to the expected harm due to his use of the substandard tank, \$1 million.

Several factors are relevant to the choice between act-based and harm-based sanctions. First, act-based sanctions need not be as high to accomplish a given level of deterrence, and thus offer an advantage over harm-based sanctions because of limitations in parties' assets. In the example in the preceding paragraph, the owner of the storage tank might be able to pay the \$1 million required if sanctions are act-based (assuming for simplicity that individuals are always found liable) but not the \$10 million required if sanctions are harm-based. Second, and closely related, because act-based sanctions need not be as high to accomplish deterrence, they tend to be preferable to harm-based sanctions when parties are risk averse. Third, act-based sanctions and harm-based sanctions may differ in the ease with which they can be applied. In some circumstances, act-based sanctions may be simpler to impose (it might be less difficult to determine whether an oil shipper properly maintains its vessels' holding tanks than to detect whether one of the vessels leaked oil into the ocean); in other circumstances, harm-based sanctions may be more readily applied (a driver who causes an accident might be caught more

easily than one who speeds but does not cause an accident). Fourth, it may be hard to calculate the expected harm due to an act, but relatively easy to ascertain the actual harm if it eventuates; if so, this constitutes an advantage of harm-based liability.³⁶

13. PRECAUTIONS

In this section we consider a model in which harm is accidental, as in the previous section, and in which the probability of harm depends on the level of precautions taken by a potential injurer. Thus, the major difference from the basic model considered in earlier sections is that the act is continuously variable. The main results of the basic analysis carry over to the model of precautions. For simplicity, we focus on the case in which enforcement is certain. Let

x = level of precautions taken by a potential injurer; and

$q(x)$ = probability of harm given x ; $q'(x) < 0$; $q''(x) > 0$.

The usual social objective of maximizing social welfare now can be expressed as minimizing social costs, that is, minimizing the sum of the cost of precautions and the expected harm: $x + q(x)h$. Let $x^* > 0$ be the solution to this problem.

First consider strict liability when the sanction is a fine equal to the harm. Then an individual's problem is to minimize $x + q(x)h$, so he will choose x^* (and obviously would not if the fine did not equal harm).

Next consider fault-based liability when the standard corresponds to x^* and the sanction is a fine equal to harm. If an individual takes less precaution than x^* , he bears costs of $x + q(x)h$, while if he takes precaution equal to or greater than x^* , he bears cost of x . It is straightforward to

³⁶ Act-based versus harm-based enforcement is discussed in Shavell (1993).

show that he will exercise precautions equal to x^* .³⁷

Thus, as in the basic theory, strict liability and fault-based liability result in the first-best outcome when sanctions are monetary and are applied for sure. Similar reasoning would demonstrate that all of the other primary conclusions in the basic theory would carry over to the model of precautions.

To illustrate, reconsider the case when individuals are risk neutral, liability is strict, and the sanction is a fine. Social costs are $x + q(x)h + e$. The level of precautions x is determined by the individual minimizing $x + q(x)p(e)f$. Again, the optimal fine must be the individual's wealth w (otherwise, f could be raised and e lowered without affecting deterrence) and the optimal p is such that $pw < h$.³⁸ Therefore, at the optimum, the level of precautions is less than the first-best level.

14. ACTIVITY LEVEL

We have been assuming that the sole decision that an individual makes is whether to act in a way that causes harm when engaging in some activity. In many contexts, however, an individual also makes a choice about his *activity level* — that is, not only does he choose whether to act in a harmful way while engaging in an activity, he also chooses whether to engage in that activity, or, more generally, at what level to do so. For example, in addition to deciding whether to comply with auto emissions controls (maintaining a catalytic converter), an individual also

³⁷ Conditional on choosing $x \geq x^*$, his best choice clearly is x^* , as that minimizes his expense, which is x^* . If he chooses $x < x^*$, his expense is $x + q(x)h$, which exceeds $x^* + q(x^*)h > x^*$. Hence, x^* is strictly optimal for him.

³⁸ Let $x(e)$ be the x determined by the individual's optimization problem, given $f = w$ and enforcement expenditures e . The individual minimizes $x + q(x)p(e)w$, with the resulting first-order condition $1 + q'(x)p(e)w = 0$. The social problem is to minimize $x(e) + q(x(e))h + e$. The first-order condition can be written as $x'(1 + q'h) + 1 = 0$. We know that $x' > 0$. Therefore, it must be that $1 + q'h < 0$. Solving for q' from the individual's first-order condition and substituting it into this expression implies that $pw < h$.

chooses how many miles to drive; the number of miles driven is the individual's level of activity.

Similarly, not only does a firm decide whether to comply with workplace safety regulations, it also chooses its level of production; the output of the firm is its level of activity.

The socially optimal activity level is such that the individual's marginal utility from the activity just equals the marginal expected harm caused by the activity. Thus, the optimal number of miles driven is the level at which the marginal utility of driving an extra mile just equals the marginal expected harm per mile driven. The determination of the optimal level of activity presumes that individuals act optimally when engaging in the activity — for example, that they drive with appropriate care.

To illustrate this formally, let

r = level of activity; and

$U(r)$ = utility from activity level r ; $U'(r) > 0$; $U''(r) < 0$.

We suppose that an individual chooses how much precaution to take (see the previous section) while engaging in a harm-creating activity, with the level of harm being proportional to his level of activity.³⁹ Then social welfare is

$$U(r) - r[x + q(x)h]. \quad (20)$$

Note that the optimal level of precaution x^* minimizes $x + q(x)h$, and thus is as discussed in section 13. The optimal level of activity therefore is determined by

$$U'(r) = x^* + q(x^*)h; \quad (21)$$

that is, the marginal utility from the activity equals the social cost of the activity, which is the sum of the cost of precautions and the expected harm.

³⁹ We are employing the model of precautions described in the previous section for convenience. It also would be possible to develop the points about activity level using the model from the basic analysis, in which the harm-producing action of an individual is not continuously variable.

Will parties' choices about their activity levels and precautions be socially correct under the two standards for imposing sanctions? The answer is that under strict liability, their choices about both activity levels and precautions will be socially correct. This is clear since, assuming for simplicity that enforcement is certain and the sanction is a fine equal to harm, their objective also is to maximize (20). In particular, because they bear a fine equal to harm, they choose the optimal level of precautions x^* when engaging in their activity. And since they incur the full social costs of precautions plus expected harm when engaging in their activity, they choose the optimal level of activity.

Under the fault-based standard, however, parties will participate in activities to a socially excessive extent. To explain, observe that parties choose the optimal level of precautions x^* in order to avoid the fine, as seen in section 13. Because parties choose this level of precautions, they will not be found liable for having violated the standard if harm occurs. Hence, their choice of activity level r is determined by maximizing $U(r) - rx^*$, with the corresponding first-order condition

$$U'(r) = x^*. \quad (22)$$

Comparing this to (21), it is evident that r exceeds r^* . The reason is that the private marginal cost of increasing participation in the activity is only the precaution cost x^* ; it does not include the expected harm. Thus, for instance, if a person complies with auto emissions standards, he will not be concerned with the fact that the more he drives, the more pollution he causes (assuming that some pollution occurs even if one obeys emissions standards). Consequently, he will drive too much.

The implication of the preceding points in relation to firms is that under the strict sanctioning rule, the product price will reflect both the cost of precautions and the expected harm

caused by production, so that the price will include the full social cost of production. Hence, the amount purchased, and thus the level of production, will be socially optimal. Under the fault-based rule, however, the product price will reflect the cost of precautions but not expected harm; thus, the amount sold, and the level of production, will be excessive.⁴⁰

These conclusions about firms are of widespread applicability. Notably, safety regulations and other regulatory requirements are often framed as standards of care that have to be met, but which, if met, free the regulated party from penalties. Hence, regulations of this character are subject to the criticism that they lead to excessive levels of the regulated activity. Making firms strictly liable for harm would be superior to safety regulation with respect to inducing socially correct activity levels.

That parties choose an excessive level of activity under the fault-based rule — of which regulation is one variant — but not under the strict liability rule, constitutes a fundamental advantage of the latter rule. This advantage is stronger the greater is the harm from engaging in the activity (given that precautions are optimal when engaging in the activity). Thus, for activities for which expected harm is likely to be substantial, the disadvantage of the fault-based standard is significant.

More generally, the advantage of strict liability over fault-based liability applies to any dimension of behavior that affects expected harm but that is not included in the definition of fault. For example, suppose that pollution damage depends both on whether a scrubber is installed as well as on the degree of care with which it is cleaned. Because the existence of a scrubber is easy to verify but its maintenance might not be, a fault-based sanctioning system

⁴⁰ Our discussion here about activity-level considerations in the context of public enforcement closely parallels the analysis of activity-level issues in the context of tort liability. See generally Shavell (1980) and Polinsky (1980b).

might, of necessity, reflect only the first dimension of behavior. Consequently, the firm will have a socially inadequate incentive to clean the scrubber under a fault standard. This problem does not arise under a strict sanctioning system because the firm has to pay for harm regardless of its cause.

15. ERRORS

Errors of the two classic types can occur in public enforcement of law. First, an individual who should be found liable might mistakenly not be found liable — what we will refer to as “mistaken acquittal.” Second, an individual who should not be found liable might mistakenly be found liable — “mistaken conviction.” For an individual who has been detected, let

ε_A = the probability of mistaken acquittal; and

ε_C = the probability of mistaken conviction.

For example, suppose police randomly monitor drivers by stopping them and administering a blood-alcohol test. The test might understate the amount of alcohol in the driver’s blood and result in mistaken acquittal, or overstate the amount and lead to mistaken conviction.

We initially consider the effect of mistake in the basic model of enforcement, assuming that the sanctioning standard is strict, the sanction is a fine, and individuals are risk neutral. Given the probability of detection p and the chances of mistaken acquittal and conviction, an individual will commit the wrongful act if and only if his gain net of his expected fine if he does commit it exceeds what he bears if he does not commit it:

$$g - p(1 - \varepsilon_A)f > -p\varepsilon_C f, \quad (23)$$

or, equivalently, if and only if⁴¹

$$g > (1 - \varepsilon_A - \varepsilon_C)pf. \quad (24)$$

Note initially that both types of error reduce deterrence: the right-hand side of (24) is declining in both ε_A and ε_C . Mistaken acquittal diminishes deterrence because it lowers the expected fine if an individual violates the law. Mistaken conviction also lowers deterrence because it reduces the difference between the expected fine from violating the law and not violating it. In other words, the greater is ε_C , the smaller the increase in the expected fine if one violates the law, making a violation less costly to the individual.⁴²

Because mistakes dilute deterrence, they tend to reduce social welfare. Specifically, to achieve any level of deterrence, it may be necessary to raise the probability of detection or the magnitude of a costly sanction to offset the effect of errors.

Now consider the optimal choice of the fine. If the probability of detection is assumed to be fixed, the dilution in deterrence caused by errors requires a higher fine to restore deterrence, so the optimal fine is higher.⁴³ If both the probability and the fine are policy instruments, the optimal fine remains maximal despite mistakes. The explanation is essentially that given previously: If the fine f were less than maximal, then f could be raised and the probability p lowered so as to keep deterrence constant, but saving enforcement costs.

If individuals are risk averse, however, the possibility of mistakes does affect the optimal fine. As we emphasized in section 7.2, the optimal fine generally is less than maximal when individuals are risk averse — lowering the fine reduces the bearing of risk. Introducing the

⁴¹ We assume that $1 - \varepsilon_A - \varepsilon_C > 0$, so that the probability that a guilty person will be found liable, $1 - \varepsilon_A$, exceeds the probability that an innocent person will be found liable, ε_C .

⁴² This point was first emphasized by Png (1986).

⁴³ Specifically, to achieve first-best behavior, it must be that $(1 - \varepsilon_A - \varepsilon_C)pf = h$, which implies that f must be higher the greater are either of the errors, ε_A or ε_C .

possibility of mistakes may increase the desirability of lowering the fine because, due to mistaken conviction, individuals who do not violate the law are subject to the risk of having to pay a fine. Indeed, because the number of persons who do not violate the law often would far exceed the number who do, the desire to avoid imposing risk on the former group can lead to a substantial reduction in the optimal fine.⁴⁴

The possibility of mistakes generally affects the optimal probability of detection. On one hand, the deterrence-diluting effects of mistakes means, as we noted, that a higher probability of detection may be needed to achieve any given level of deterrence; this effect tends to raise the optimal expenditure on enforcement. On the other hand, because mistakes reduce the productivity of enforcement expenditures by a factor of $1 - \varepsilon_A - \varepsilon_C$ (see (24)), the cost of achieving a given level of deterrence is higher; this effect tends to reduce the optimal expenditure on enforcement. Either of these effects could dominate and lead to an optimal probability of detection that is higher or lower than in the absence of mistakes.

Next, consider imprisonment and mistake. As in the case of fines, mistakes of both type dilute the deterrent effect of imprisonment. Additionally, as in the case without mistakes, the optimal imprisonment term is maximal if individuals are risk neutral or risk averse in imprisonment, but is generally not maximal if they are risk preferring in imprisonment.⁴⁵

The possibility of mistakes also affects individuals' decisions regarding their participation in activities. Everything else equal, mistaken acquittals lead to increased

⁴⁴ Building on Polinsky and Shavell (1979), Block and Sidak (1980, pp. 1135-1139) emphasize the desirability of lowering sanctions on risk-averse individuals because of mistakes.

⁴⁵ That the optimal term remains maximal if individuals are risk neutral or risk averse might seem surprising because one might expect that the chance of mistaken conviction would result in a lower optimal term. But the usual argument still applies: If the term were not maximal, it could be raised and the probability of detection could be lowered at least proportionally without sacrificing deterrence. Hence, the aggregate amount of jail time served by individuals who do not commit the harmful act would remain the same or fall, and enforcement expenditures would decline.

engagement in the activity, and mistaken convictions result in decreased engagement. The net effect could be a socially excessive or inadequate activity level.⁴⁶

We have not yet commented on fault-based liability and mistake. In this context, an important implication of mistake is that some individuals will bear sanctions even if they comply with the fault standard. Consequently, both types of error reduce the incentive to comply with the fault standard in the basic model, for essentially the same reasons as under strict liability.

However, in the model with variable precautions, there is a possibility that error will lead to an excessive level of precautions — above the optimal level. Assume that the actual level of precautions x is observed with an error ε . If the observed level of precautions, $x + \varepsilon$, exceeds the standard \bar{x} , the person will not be liable because he will not be found to be at fault. But if $x + \varepsilon$ is less than \bar{x} , then the person will be liable. The person will be found mistakenly liable if x is greater than or equal to \bar{x} but $x + \varepsilon$ is less than \bar{x} ; the person will be mistakenly acquitted if x is less than \bar{x} but $x + \varepsilon$ equals or exceeds \bar{x} . Note that the probabilities of the two types of error depend on the person's choice of x . By choosing an x above \bar{x} , the person can reduce the risk of mistaken conviction, and it can be shown that he will be led to choose such an excessive x under fairly general conditions. In other words, individuals will often have a motive to take excessive precautions in order to reduce the chance of erroneously being found at fault.⁴⁷

Errors also influence individuals' participation in the activity under fault-based liability, similar to the effects of errors under strict liability. The main difference is that there is a general tendency for individuals to participate in the activity to an excessive extent under fault-based

⁴⁶ Recall, too, that if the probability of detection is fixed, the fine needs to be raised to offset the deterrence-diluting effects of mistakes. Raising the fine to this extent, however, leads to an inadequate incentive to engage in the activity. This problem, in turn, can be remedied by use of an appropriate subsidy for participating in the activity. For the details behind this point, see Png (1986).

⁴⁷ This point was first emphasized by Craswell and Calfee (1986).

liability, for the reason explained in the previous section. Only if the chance of mistaken conviction is sufficiently high would this conclusion be reversed.

Finally, observe that the probabilities of error can be influenced by policy choices. For example, prosecutorial resources can be increased in order to reduce the probability of mistaken acquittal, or the standard of proof can be raised to reduce the chance of mistaken conviction (although this presumably increases the likelihood of mistaken acquittal). Because the reduction of both types of error increase deterrence, expenditures made to reduce errors may be socially beneficial.⁴⁸

16. COSTS OF IMPOSING FINES

We inquire in this section about the implications of costs borne by enforcement authorities in imposing fines.⁴⁹ Our principal observation is that such costs should raise the level of the fine.

To elaborate, suppose that the probability of detection is fixed at p , that liability is strict, and that individuals are risk neutral. If fines are costless to impose, the optimal fine is h/p , the harm divided by the probability of detection (see (14)). Now suppose that the enforcement authority bears a cost each time a fine is imposed; let

k = cost of imposing a fine on an offender.

It is easy to verify that the optimal fine then is

$$f^* = h/p + k; \tag{25}$$

the cost k should be added to the fine that would otherwise be desirable. The explanation is that,

⁴⁸ On the value of accuracy in adjudication, see Kaplow and Shavell (1994a).

⁴⁹ We have already discussed the cost of imposing imprisonment sanctions — specifically, the cost to the state per unit of the imprisonment term, c .

if an individual commits a harmful act, he causes society to bear not only the immediate harm h , but also, with probability p , the cost k of imposing the fine — that is, his act results in an expected total social cost of $h + pk$. If the fine is set according to (25), the individual's expected fine is $h + pk$, which leads him to engage in the harmful act if and only if his gain exceeds the expected total social cost of the act.

There may be other costs associated with the imposition of fines. In particular, suppose that detection is followed by a costly second stage during which the state investigates and prosecutes an individual, and at the end of which a fine is imposed only with a probability. Let

s = cost of the investigation-prosecution stage; and

q = probability of a fine being imposed after the investigation-prosecution stage.

Hence, the probability that an individual will have to pay a fine is pq and the expected costs of imposing a fine, including the expected investigation-prosecution cost, become $ps + pqk$.

It is readily shown that the optimal fine now is

$$f^* = h/pq + s/q + k. \quad (26)$$

This formula illustrates a general principle: the optimal fine equals the costs incurred by society as a result of the harmful act divided by the probability — at the time that each component of cost is incurred — that the individual will have to pay the fine. Thus, h is divided by pq because, when the harm occurs, the probability of having to pay the fine is pq ; and s is divided by q because, when the investigation-prosecution costs are incurred, the probability of having to pay the fine is q . If the fine is computed according to this principle, the expected fine will equal the expected social costs due to an individual committing a harmful act, including the harm caused and the expected sanctioning costs — that is, $h + ps + pqk$.

Note that under fault-based liability, the costs of imposing fines is significantly lower, if not zero. This is because, if individuals comply with the fault standard, they do not bear

sanctions, in which case there are no costs associated with imposing sanctions. However, if individuals are found at fault (say because of errors), the fines imposed on them also should reflect the costs of imposing fines.

Finally, observe that not only does the state incur costs when fines are imposed, so do the individuals who pay the fines (such as legal defense expenses). The costs borne by individuals, however, do not affect the formula for the optimal fine. Individuals properly take these costs into account because they bear them directly.⁵⁰

17. GENERAL ENFORCEMENT

In many settings, enforcement may be said to be *general* in the sense that several different types of violations may be detected by an enforcement agent's activity. For example, a police officer waiting along the side of a road may notice a driver who litters as well as a driver who goes through a red light or who speeds; or a tax auditor may detect a variety of infractions when he examines a tax return. To analyze this type of situation, suppose that a single probability of detection applies to all harmful acts, regardless of the magnitude of the harm.⁵¹ (The contrasting assumption is that enforcement is *specific*, meaning that the probability is chosen independently for each type of harmful act.)

The main point that we want to make is that when enforcement is general, the optimal sanction rises with the severity of the harm and is maximal only for relatively high harms. To see this, assume that liability is strict, the sanction is a fine, and individuals are risk neutral. Let

⁵⁰ The points developed in this section were first presented in Polinsky and Shavell (1992), although early writers on enforcement theory — including Becker (1968, p. 192) and Stigler (1970, p. 533) — recognized that sanctions should reflect enforcement costs.

⁵¹ It will be clear that the main point developed in this section does not depend on the assumption that the same probability applies to all acts. The only requirement is that the probabilities for different acts are linked, all a function of the same enforcement expenditure.

$f(h)$ be the fine given harm h . Then, for any general probability of detection p , the optimal fine schedule is

$$f^*(h) = h/p, \tag{27}$$

provided that h/p does not exceed the maximum feasible fine (say individuals' wealth level w); if h/p is not feasible, the optimal fine is maximal. This schedule is obviously optimal given p because it implies that the expected fine equals harm, thereby inducing first-best behavior, whenever that is possible.

The question remains whether it would be desirable to lower p and raise fines to the maximal level for the low-harm acts for which $f^*(h)$ is less than maximal. But if p is reduced for the relatively low-harm acts (and the fine raised for them), then p — being general — is also reduced for the high-harm acts for which the fine is already maximal, resulting in lower deterrence of these acts. The decline in deterrence of high-harm acts may cause a greater social loss than the savings in enforcement costs from lowering p . To express this point differently, p must be sufficiently high to avoid significant underdeterrence of high-harm acts (for which fines are maximal). But since this p also applies to less harmful acts, the fines for them do not need to be maximal in order to deter them appropriately.⁵²

The result that, when enforcement is general, sanctions should rise with the severity of harm up to a maximum also holds if the sanction is imprisonment and if liability is fault-based. The underlying reasoning is the same as that given above.⁵³

⁵² Note that if p could be varied independently for a low-harm act and for a high-harm act — that is, if enforcement is specific rather than general — then it would be desirable to lower p and raise the fine for a low-harm act if the fine for it were less than maximal.

⁵³ The basic point of this section was first made by Shavell (1991b); see also Mookherjee and Png (1992) for a closely related analysis.

18. MARGINAL DETERRENCE

In many circumstances, an individual may consider which of *several* harmful acts to commit, for example, whether to release only a small amount of a pollutant into a river or a large amount, or whether only to kidnap a person or also to kill him. In such contexts, the threat of sanctions plays a role in addition to the usual one of deterring individuals from committing harmful acts: for individuals who are not deterred, expected sanctions influence *which* harmful acts individuals choose to commit. Notably, such individuals will have a reason to commit less harmful rather than more harmful acts if expected sanctions rise with harm. Deterring a more harmful act by having its expected sanction exceed that for a less harmful act is sometimes referred to as *marginal deterrence*.⁵⁴

Other things being equal, it is socially desirable that enforcement policy creates marginal deterrence, so that those who are not deterred from committing harmful acts have a reason to moderate the amount of harm that they cause. This suggests that sanctions should rise with the magnitude of harm and, therefore, that most sanctions should be less than maximal. However, promoting marginal deterrence may conflict with achieving deterrence generally: for the schedule of sanctions to rise steeply enough to accomplish marginal deterrence, sanctions for less harmful acts might have to be so low that individuals are inadequately deterred from committing these acts.⁵⁵

To illustrate the implications of marginal deterrence, consider the following example in which sanctions are monetary and liability is strict. Suppose that there are two harmful acts, with

⁵⁴ The notion of marginal deterrence was remarked upon in some of the earliest writing on enforcement; see Beccaria (1767, p. 32) and Bentham (1789, p. 171). The term *marginal deterrence* apparently was first used by Stigler (1970).

⁵⁵ For formal treatments of marginal deterrence, see Shavell (1992), Wilde (1992), and Mookherjee and Png (1994).

harms h_1 and h_2 , where $h_1 < h_2$, that the probability of detection p is the same for both acts, and that individuals have the same level of wealth w . We will first consider a one-act model in which there can be no marginal deterrence because each individual can commit only one type of harmful act. We will then compare the results in this case to a two-act model in which each individual can commit either of two harmful acts.

In the one-act model, suppose some individuals have the opportunity to commit an act causing harm of h_1 and other individuals have the opportunity to commit an act causing harm of h_2 . It is optimal to set the fine for the high-harm act equal to w , for otherwise it would be possible to raise the fine for both acts and lower the probability of detection p without affecting deterrence, but saving enforcement costs. It also follows that the optimal p is such that pw is less than h_2 , that is, there is some underdeterrence of the high-harm act. The reason is that if $pw = h_2$, there would be no first-order loss of social welfare in terms of gains and harm if p is lowered (since marginal individuals are those for whom $g = h_2$), but enforcement costs would be saved. Given the common probability p , the fine f_1 for the lesser offense then can be set such that $pf_1 = h_1$ (assuming such a fine is feasible), achieving first-best deterrence of this offense.

In the two-act model, each individual can commit an act causing harm of h_1 or an act causing harm of h_2 . Again, it is optimal to set the fine for the high-harm act equal to w and for there to be underdeterrence of the high-harm act. Now, however, f_1 should be such that pf_1 is less than h_1 , instead of equaling h_1 . The essential reason for this result is that the reduction in f_1 from h_1/p leads some offenders to commit the act causing harm h_1 instead of the act causing higher harm h_2 . This can be shown to raise social welfare even though the reduction in f_1 leads some individuals to commit the low-harm act who otherwise would not have committed either

harmful act.⁵⁶ In other words, achieving marginal deterrence by reducing the expected fine for the low-harm act raises social welfare.

Two additional observations should be made about marginal deterrence. First, marginal deterrence can be promoted by increasing the probability of detection as well as the magnitude of sanctions for acts that cause greater harm. For example, kidnappers can be deterred more from killing their victims if greater police resources are devoted to apprehending kidnappers who murder their victims than to those who do not. (Note, though, that in circumstances in which enforcement is general, the probability of detection cannot be independently altered for acts that cause different degrees of harm.) Second, marginal deterrence is naturally accomplished if the expected sanction equals harm for all levels of harm; for if a person is paying for harm done, he will have to pay appropriately more if he does greater harm. Thus, for instance, if a polluter's expected fine would rise from \$100 to \$500 if he dumps five gallons instead of one gallon of waste into a lake, where each gallon causes \$100 of harm, his marginal incentives to pollute will be correct.⁵⁷

19. PRINCIPAL-AGENT RELATIONSHIP

Although we have assumed that an offender is an independent single actor, in fact the offender is often an agent of a principal. For example, the agent could be an employee of a firm

⁵⁶ This conclusion essentially follows from two observations. First, because pw is less than h_2 , some individuals who had been committing the high-harm act were causing a net loss of social welfare (their gain was less than h_2). Second, as f_l is lowered marginally from h_l/p , individuals who are induced to commit the act causing harm h_l (who either would have committed the high-harm act or not committed any harmful act) cause no net loss of social welfare (their gain equaled h_l).

⁵⁷ As we discussed in section 7.1, however, it generally is desirable for society to tolerate some underdeterrence in order to save enforcement costs, in which case expected sanctions will be less than harm. Then, consideration of marginal deterrence alters the structure of sanctions that would otherwise be best, as the comparison of the one-act model to the two-act model in this section showed.

or a subcontractor of a contractor. The enforcement problem is now how to maximize social welfare by choosing enforcement effort and the sanctions to be imposed on principals and agents. This maximization is carried out under the assumption that a principal chooses a contract with his agent that maximizes the principal's expected utility subject to two constraints: that the agent receive his reservation level of expected utility; and that the agent maximizes his own expected utility.

When harm is caused by an agent, many of our conclusions from the basic analysis carry over if the sanction is imposed on the principal. For example, given the probability of detection p , it is optimal for a risk-neutral principal to face a fine of h/p . Then the expected fine is equal to harm done. Consequently, the principal will behave socially optimally in controlling his agents, and in particular will contract with them and monitor them in ways that will give the agents socially appropriate incentives to reduce harm.⁵⁸

A question about enforcement that arises when there are principals and agents is how to allocate financial sanctions between them. First observe that the particular allocation of sanctions may not matter when, as would be the natural presumption, the principal and the agent can reallocate sanctions through their own contract. For example, if the agent finds that he faces a large fine but is more risk averse than the principal, the principal can assume it; conversely, if the fine would be imposed on the principal, he can bear that risk and not impose an internal sanction on the agent. Thus, the post-contract penalties that the agent suffers may not be affected

⁵⁸ There is relatively little literature on the question of optimal enforcement when wrongdoers are agents of principals. Newman and Wright (1990) study the optimal monetary sanction to impose on a risk-neutral principal when liability is strict and is imposed for sure; they show that it equals harm. Polinsky and Shavell (1993) demonstrate the potential desirability of imposing criminal sanctions on an employee who causes harm, even when the employer is capable of paying for the harm. Arlen (1994) examines the effect of sanctions on corporations' incentives to monitor their employees, and she emphasizes the possibility that corporations may have perverse incentives not to monitor if they would become liable as a result of their discovering and reporting employee violations. Also, Shavell (1997b) finds that optimal sanctions on corporations could be above or below harm when employee assets are less than harm.

by the particular division of sanctions initially selected by the enforcement authority.

The allocation of monetary sanctions between principals and agents would matter, however, if some allocations allow the pair to reduce their total burden. An important example is when a fine is imposed only on the agent and he is unable to pay it because his assets are less than the fine.⁵⁹ Then he and the principal (who often would have higher assets) would jointly escape part of the fine, diluting deterrence. The fine therefore should be imposed on the principal rather than on the agent (or at least the part of the fine that the agent cannot pay).

A closely related point is that the imposition of imprisonment sanctions on agents may be desirable when their assets are less than the optimal fine, even if the principal's assets are sufficient to pay the fine. The fact that an agent's assets are limited means that the principal may be unable to control him adequately through the use of contractually-determined penalties, which can only be monetary. For example, a firm may not be able, despite the threat of salary reduction or dismissal, to induce its employees never to rig bids. In such circumstances, it may be socially valuable to use the threat of personal criminal liability and a jail sentence to better control agents' misconduct.⁶⁰

20. SETTLEMENTS

We have thus far assumed that when an individual who should be found liable is discovered, he will be sanctioned in some automatic fashion. In practice, however, an individual must be found liable in a trial, and before this occurs, it is common for an individual to settle in lieu of trial. (In the criminal context, the settlement usually takes the form of a *plea bargain*, an

⁵⁹ See Sykes (1981) and Kornhauser (1982).

⁶⁰ This point is discussed by Segerson and Tietenberg (1992) and emphasized by Polinsky and Shavell (1993).

agreement in which the individual pleads guilty to a reduced charge.) Given the prevalence of settlements, it is important to consider how they affect deterrence and the optimal system of public enforcement, and whether settlements are socially desirable.

A general reason why a wrongdoer who has been caught might prefer an out-of-court settlement to a trial is that a settlement saves him time and/or money. Public enforcers presumably would prefer settlements for this reason as well. To amplify, consider a risk-neutral detected individual subject to a fine f and let

r = probability of conviction;

c_I = individual's litigation costs; and

c_P = prosecutor's litigation costs.

Assume that the parties agree on the probability of conviction. If the case goes to trial, the expected cost to the individual is $rf + c_I$, and the expected gain to the prosecutor, assuming his goal is to maximize expected penalties imposed net of his litigation costs, is $rf - c_P$. Thus, any settlement resulting in a fine between $rf - c_P$ and $rf + c_I$ would make both parties better off because the cost of litigation would be avoided. The same point would apply if the individual were subject to a jail sentence rather than a fine. Note, however, that if the parties disagree about the probability of conviction, a settlement might not occur — specifically, if the individual is relatively optimistic and believes that the probability of his being convicted is sufficiently less than the prosecutor thinks it is.

A second benefit of a settlement is that it eliminates the risks inherent in the trial outcome, a benefit to parties who are averse to such risks.⁶¹

⁶¹ These benefits of settlement are well-recognized in the economic literature on civil litigation; see the surveys by Cooter and Rubinfeld (1989) and Spier (2006). For early discussions of settlement in the context of public enforcement, see Landes (1971) and Grossman and Katz (1983), and more recently, see, for example, Reinganum (1988), Polinsky and Rubinfeld (1989), Kobayashi and Lott (1992), and Miceli (1996).

The preceding advantages of settlement to the parties suggest that settlement is socially valuable, but the effect of settlement on deterrence is a complicating factor. Specifically, settlements dilute deterrence: for if individuals desire to settle, it must be because the expected disutility of sanctions is lowered for them. However, because settlements reflect the sanctions that would be imposed at trial, the state may be able to offset this settlement-related reduction in deterrence by increasing the level of sanctions. If so, settlements need not compromise the overall level of deterrence.⁶²

Settlements may have other socially undesirable consequences. First, they may result in sanctions that are not as well tailored to harmful acts as would be true of court-determined sanctions. Second, settlements hinder the amplification and development of the law through the setting of precedents. Third, settlements also sometimes allow individuals to keep aspects of their behavior secret, which can reduce deterrence. Fourth, settlements for prison terms can result in increases in public expenditures on jail if individuals are risk averse in imprisonment.⁶³ A prosecutor whose goal is to maximize social welfare, as opposed to maximizing the expected sanction less prosecution costs, presumably would take these additional factors into account and sometimes refuse to settle even though the settlement saves litigation costs and avoids risk.⁶⁴

⁶² The deterrence-diluting effects of settlement and other aspects of the social desirability of settlement have been discussed in the private litigation context by Polinsky and Rubinfeld (1988), Shavell (1997a), and Spier (1997). A related discussion in the public enforcement context appears in Polinsky and Rubinfeld (1989).

⁶³ For example, suppose a defendant faces a 50% chance of a 5 year sentence and a 50% chance of a 15 year sentence, with an expected sentence of 10 years. If he is risk averse, he will strictly prefer a certain sentence of 10 years. This implies that if prosecutors want to maintain deterrence, they must demand a settlement of more than 10 years, say 12 years, which increases the cost of imprisonment.

⁶⁴ The question of what prosecutors maximize has received almost no attention from law and economics scholars, although two exceptions are Miceli (1996) and Glaeser et al. (2000).

21. SELF-REPORTING

We have assumed that individuals are subject to sanctions only if they are detected by an enforcement agent, but in fact parties sometimes disclose their own violations to enforcement authorities. For example, firms often report violations of environmental and safety regulations, individuals usually notify police of their involvement in traffic accidents, and even criminals occasionally turn themselves in.

Self-reporting can be induced by lowering the sanction for individuals who disclose their own infractions. To avoid significantly reducing deterrence, however, the reward for self-reporting can be made relatively small. For example, suppose that the fine if an individual does not self report is \$1,000 and that the probability of detection is 10%, so the expected fine is \$100. If the fine if one self-reports is \$99, individuals will self-report but deterrence will barely be reduced.

To express this formally, assume for simplicity that individuals are risk neutral, and suppose that if an individual commits a violation and does not self-report, his expected fine is pf .

Let

f' = fine if a violator self-reports,

and set

$$f' = pf - \varepsilon, \quad (28)$$

where $\varepsilon > 0$ is small. A violator will therefore want to self-report because f' is less than pf , but the deterrent effect of the sanction will approximate that if he did not self-report.

Given that self-reporting can be induced essentially without compromising deterrence, why is self-reporting socially advantageous? One reason is that self-reporting lowers enforcement costs because the enforcement authority does not have to identify and prove who

the violator was. Environmental enforcers do not need to spend as much effort trying to detect pollution and establishing its source if firms that pollute report that fact.⁶⁵ Second, self-reporting reduces risk, and thus is advantageous if individuals are risk averse. Drivers bear less risk because they know that if they cause an accident, they will be led to report this to the police and suffer a lower and certain sanction (of approximately pf), rather than face a substantially higher sanction imposed only with some probability. Third, self-reporting may allow harm to be mitigated. Early identification of a toxic leak, for example, will facilitate its containment and clean-up.⁶⁶

22. REPEAT OFFENDERS

In practice, the law often sanctions repeat offenders more severely than first-time offenders. For example, under the U.S. Sentencing Commission's sentencing guidelines for Federal crimes, both imprisonment terms and criminal fines are enhanced if an offender has a prior record. Civil money penalties also sometimes depend on whether the offender has a record of prior offenses. We explain here why such policies may be socially desirable.

Note first that sanctioning repeat offenders more severely cannot be socially advantageous if deterrence always induces first-best behavior. If the sanction for polluting and causing a \$1,000 harm is \$1,000, then any person who pollutes and pays \$1,000 is a person

⁶⁵ In some contexts, however, self-reporting will not save enforcement costs. For example, suppose that a police officer waits by the side of a road to spot speeders. Then, were a driver to report that he had sped, this would not reduce policing costs, presuming that the officer still needs to be stationed at the roadside to watch for other speeders. Usually, though, there would be some cost savings as a result of self-reporting (for example, the police officer would not have to chase as many speeders).

⁶⁶ The basic theory of self-reporting in public enforcement is developed in Kaplow and Shavell (1994b); see also Malik (1993) and Innes (1999). Related literature concerns the reporting of income by individuals to tax authorities and the reporting of costs by regulated firms to regulatory authorities. See, for example, Andreoni et al. (1998) and Laffont and Tirole (1993).

whose gain from polluting (say the savings from not installing pollution control equipment) must have exceeded \$1,000. Social welfare therefore is higher as a result of his polluting. If such an individual polluted and was sanctioned in the past, that only means that it was socially desirable for him to have polluted previously. Raising the sanction because of his having a record of prior convictions would overdeter him now.

Accordingly, only if deterrence is inadequate is it possibly desirable to condition sanctions on offense history to increase deterrence. But deterrence often will be inadequate because, as we emphasized in section 7.1, it will usually be worthwhile for the state to tolerate some underdeterrence in order to reduce enforcement expenses.

Given that there is underdeterrence, making sanctions depend on offense history may be beneficial for two reasons. First, the use of offense history may create an additional incentive not to violate the law: if detection of a violation implies not only an immediate sanction, but also a higher sanction for a future violation, an individual will be deterred more from committing a violation presently.⁶⁷ Second, making sanctions depend on offense history allows society to take advantage of information about the dangerousness of individuals and the need to deter them: individuals with offense histories may be more likely than average to commit future violations, which might make it desirable for purposes of deterrence to impose higher sanctions on them.⁶⁸

⁶⁷ There is a subtlety in demonstrating the optimality of punishing repeat offenses more severely. Namely, if there is a problem of underdeterrence, one might wonder why it would not be optimal to raise the sanction to the maximum level for every offense (in which case repeat offenses would not be punished more severely). It must be shown that punishing all offenses maximally is inferior to punishing first offenses less than maximally and punishing repeat offenses more severely. See Polinsky and Shavell (1998) on the possible optimality of making sanctions depend on offense history because of the additional deterrence that such a policy creates. Miceli and Bucci (2005) supply a different reason for raising the fine for the second offense — that there is little additional social stigma associated with a second offense, so that a higher sanction is needed to maintain deterrence. Emons (2003), however, raises the possibility that it may be optimal to lower the sanction for the second offense.

⁶⁸ Note that this reason for making sanctions depend on offense history is different from the first reason: the second reason involves the assumption that offenders are different and that the optimal sanction for some offenders is higher than for others; the first reason applies even if individuals are identical. On the second, information-based, reason for making sanctions depend on offense history, see Rubinstein (1979), Polinsky and Rubinfeld (1991), and

There is also an incapacitation-based reason for making sanctions depend on offense history. Repeat offenders are more likely to have higher propensities to commit violations in the future and thus more likely to be worth incapacitating by imprisonment.⁶⁹

23. IMPERFECT KNOWLEDGE ABOUT THE PROBABILITY AND MAGNITUDE OF SANCTIONS

Although we have made the simplifying assumption that individuals know the probability of detection and the magnitude of sanctions, it is obvious that individuals frequently have imperfect knowledge of these variables. They generally possess only subjective probability distributions of the probability of a sanction and its magnitude. They might not know the true probability of a sanction for several reasons: because the enforcement authority refrains from publishing information about the probability (perhaps hoping that individuals will believe it to be higher than it is); because the probability depends on factors that individuals do not fully understand (the probability of a tax audit, for example, is influenced by factors that are kept secret from taxpayers); and because probabilities might be difficult for individuals to assess.⁷⁰ Also, individuals may have incomplete knowledge of the true magnitudes of sanctions, particularly if sanctions are not fixed by law, but are to some degree discretionary.⁷¹

The implications of individuals' imperfect knowledge are straightforward to ascertain.

Chu et al. (2000).

⁶⁹ See section 25 for a discussion of the incapacitation rationale for the use of imprisonment sanctions.

⁷⁰ On the problems that individuals have in evaluating and using probabilities, see, for example, Kahneman et al. (1982).

⁷¹ In addition, individuals could have imperfect information about the prevailing standard of liability, not being sure whether it is strict or fault-based. This type of mistake, about a discrete issue, seems less likely to be significant than errors in assessing the probability and magnitude of sanctions.

First, to predict how individuals behave, what is relevant, of course, is not the actual probability and magnitude of a sanction, but perceptions of them.

Second, to determine the optimal probability and magnitude of sanctions, account must be taken of the relationship between the actual and the perceived values. To illustrate, suppose that the perceived probability is a single value $\hat{p}(p)$, where \hat{p} is increasing in the true probability p , and, similarly, that the perceived fine is $\hat{f}(f)$, where \hat{f} is increasing in the true fine f . Thus, if the probability of detection p is fixed, the optimal fine is set such that $\hat{p}(p)\hat{f}(f) = h$, assuming such a fine is feasible. This might imply a higher or a lower fine than when perceptions are accurate.

Third, the result that the optimal fine should be maximal when individuals are risk neutral continues to hold. By raising the fine to the maximum, the perceived fine also will be maximal. The probability of detection then can be lowered, thereby saving enforcement costs without reducing deterrence.

Several other observations are worth making. One concerns lags in learning about changes in enforcement policy. For example, suppose that there is a delay of at least a year before individuals fully comprehend a change in the probability of enforcement. Then if enforcement resources are increased so as to make the probability, say, 15% rather than 10%, there might not be a significant increase in deterrence for some time, making such an investment less worthwhile.⁷² Another observation involves the difficulty in learning about variations in enforcement policy when enforcement policy is described by a distribution. For instance, suppose that the sanction for some act, such as robbery, can vary (say from one month of jail time to 10 years), and that individuals' perceptions are quite rough, not based on true averages,

⁷² Similarly, suppose that individuals treat all probabilities of enforcement that are low, say below 1%, as if they were probabilities of 1%, because it is not possible for individuals to make discriminations finer than 1%. Then if the actual probability is ½%, spending more on enforcement to make the probability 1% would not be beneficial because deterrence would not increase.

but mainly on the possible range of sanctions. Then increasing the average sentence within this range might have very little effect on deterrence. The processes by which individuals formulate probabilities of sanctions and their magnitudes are important, therefore, to determining optimal deterrence policy.⁷³

24. CORRUPTION

In this section we examine the possible corruption of law enforcement agents, how corruption lessens deterrence and distorts participation in harm-creating activities, and what policies may be employed to combat corruption.⁷⁴ One form of corruption is bribery, in which a law enforcer accepts a payment in return for not reporting a violation (or for reducing the mandated sanction for the violation). For example, in consideration of a bribe payment, a police officer may overlook a speeding violation or a building inspector may ignore a code infraction. (For simplicity, we do not distinguish between a bribe offered by an individual and an extortion demand made by the enforcer — a payment for not turning in the individual.) A second form of corruption is framing and framing-related extortion, in which an enforcement agent may frame an innocent individual or threaten to frame him in order to extort money from him.

One reason bribery is socially undesirable is that it dilutes deterrence of violations of law. This is because bribery results in a lower payment by an individual than the sanction for the offense. To be concrete, let

⁷³ Bebchuk and Kaplow (1992) consider imperfect information about the probability of sanctions and emphasize that maximal sanctions may not be socially desirable. See also Kaplow (1990a), which takes into account learning about whether acts are subject to sanctions, and Sah (1991), which focuses on the process by which individuals form perceptions of the probability of detection.

⁷⁴ The discussion in this section is based principally on Polinsky and Shavell (2001). We do not examine the corruption of the government procurement process, such as the payment of a bribe to a government agent in order to obtain a defense contract.

λ = fraction obtained by the enforcer of the surplus from a bribe agreement.

Therefore, because the surplus from a bribe agreement is the fine f , an individual pays a bribe of λf . To the degree that λ is less than one, there is underdeterrence.⁷⁵ Similarly, bribery leads to excessive participation in the harm-creating activity.

Framing and framing-related extortion also dilute deterrence of violations of law. The reason is that framing and extortion imply that those who act innocently face an expected sanction, so that the difference between the expected sanction if individuals commit a violation and if they do not is lessened. If, for example, individuals who violate the law face an expected fine of \$1,000 and innocent individuals face an expected fine of \$200 due to the risk of being extorted or framed, the additional cost to an innocent individual of committing the offense is \$800 instead of \$1,000. (This point is essentially the same as the observation in section 15 that mistaken convictions dilute deterrence.) Additionally, framing and framing-related extortion undesirably discourage participation in the harm-creating activity.

Because corruption dilutes deterrence and distorts activity decisions, its control may be socially desirable. One way to reduce corruption is to impose fines (or imprisonment sentences) on individuals caught engaging in bribery, extortion, and framing. For example, suppose there is a fine for bribery imposed on the enforcer with a probability. Let

f_B = fine imposed on the enforcer for engaging in bribery; and

p_B = probability that an enforcer is caught engaging in bribery.

Bribery will be deterred if the surplus from a bribe agreement is eliminated, that is, if $p_B f_B \geq f$. Otherwise, the bribe payment will be $p_B f_B + \lambda(f - p_B f_B)$, which exceeds λf , so that deterrence of the offense is greater due to the sanctioning of bribery. For previously discussed reasons, the

⁷⁵ Garoupa and Klerman (2004) discuss how the threat of an imprisonment sanction for the offense will lead the offender to pay a higher bribe than otherwise, thereby reducing the deterrence-diluting effect of bribery.

optimal fine to impose on risk-neutral parties for engaging in bribery maximal, and the optimal fine for enforcers who frame innocent individuals also is maximal. But, surprisingly, framing-related extortion should not be penalized.⁷⁶

Corruption also can be reduced by paying enforcers rewards for reporting violations. Such payments will reduce their incentive to accept bribes because they will sacrifice their rewards if they fail to report violations. Indeed, sufficiently high rewards would eliminate all incentives to accept bribes. But high rewards may not be optimal because they give enforcers a greater incentive to frame innocent individuals, and high rewards tend to increase framing-related extortion payments (because enforcers sacrifice more by accepting the extortion payment). The optimal reward balances the beneficial effect of using rewards to offset the dilution of deterrence due to bribery with the detrimental effects associated with increased framing and extortion of innocent individuals.

A third way to control corruption is to pay enforcers more than their reservation wage (that is, to pay them an efficiency wage). Then they would have more to lose if punished for corrupt behavior and denied future employment. There is, however, a social cost to the state of paying enforcers more than the wage necessary to attract them — the distortions caused by the additional taxes needed to make such payments.

The discussion to this point implicitly presumed that the fine for the harmful act is fixed. A question that naturally arises, however, is whether the deterrence-diluting effects of corruption can be offset by raising the fine on offenders. For example, suppose that the optimal fine would

⁷⁶ The kernel of the reason is that penalizing framing-related extortion will lead to one of two detrimental consequences: it will either fail to deter extortion and result in higher costs to innocent individuals (the sum of their expected extortion payment and the expected fine on them for paying extortion); or else it will cause enforcers to switch from extorting money from innocent individuals to actually framing them, which is socially worse. This result is demonstrated in Polinsky and Shavell (2001), which also discusses qualifications to this point.

be \$100 if a fine were always paid when an offender is caught, but that bribery results in a bribe payment equal to \$50, one half of the fine. Could not the fine on an offender be increased to \$200, so that the bribe would then be \$100 and the effective penalty be exactly what is desired? In the basic risk-neutral model of enforcement, it is not possible to raise the fine because the optimal fine is maximal. More realistically, however, the optimal fine is less than maximal for a variety of reasons, including those related to risk aversion, marginal deterrence, and general enforcement. Then, while it would be possible to raise the fine to offset the deterrence-diluting effects of corruption, doing so would lead to social costs (for example, greater bearing of risk) and might not be desirable.⁷⁷

25. INCAPACITATION

Our discussion of public enforcement has presumed that the threat of sanctions reduces harm by discouraging individuals from causing harm — that is, by deterring them. However, a different way for society to reduce harm is by imposing sanctions that remove parties from positions in which they are able to cause harm — that is, by *incapacitating* them. Imprisonment is the primary incapacitative sanction, although there are other examples: individuals can lose their drivers licenses, preventing them from doing harm while driving; businesses can lose their right to operate in certain domains, and the like. We focus here on imprisonment, but what we say applies to incapacitative sanctions generally.

⁷⁷ Becker and Stigler (1974) focus on the control of bribery and consider paying rewards to enforcers or requiring them to post bonds that would be forfeited if they are caught engaging in bribery. Mookherjee and Png (1995) analyze bribery and conclude, given their assumption that fines are unbounded, that it is optimal to eliminate bribery. Bowles and Garoupa (1997) also discuss the control of bribery through sanctions. Hindriks et al. (1999) study bribery and extortion in the context of tax evasion, and examine rewards and penalties as methods of control. Other writing on corruption includes Pashigian (1975), Klitgaard (1988), Shleifer and Vishny (1993), Bardhan (1997), and Rose-Ackerman (1999); several of these articles focus on corruption in the awarding of government contracts and licenses rather than corruption in the imposition of sanctions for violations of law.

To better understand public enforcement when sanctions are incapacitative, suppose that their sole function is to incapacitate; that is, assume for simplicity that sanctions do not deter. (For instance, deterrence might not occur if, given the relevant range of the probabilities and magnitudes of the sanctions, individuals' gains from harmful acts exceed the expected sanctions.) Let

$h(t)$ = harm that would be caused by an individual at age t if not in jail.

We assume that $h(t)$ either is constant or declines with age.

Assuming that the social goal is to minimize the sum of the harm and the cost c of incarceration,⁷⁸ the optimal policy is to keep an individual in jail as long as $h(t) > c$. In other words, if the harm the individual would cause exceeds the cost of imprisonment, an individual should be put in prison and kept there as long as the harm continues to exceed the cost of imprisonment. He should be released otherwise. Put differently, the optimal prison term as a function of potential harm caused is zero up to a threshold — the point at which harm equals the cost of imprisonment — and then rises discontinuously to the length of time during which the person's harm if released would exceed imprisonment costs. Jail should only be used to incapacitate individuals who otherwise would have caused relatively high harm.

Two points about the incapacitative rationale are important to note. First, evidence exists suggesting that the harm caused by individuals declines with their age.⁷⁹ Thus, from the incapacitative standpoint, it often will be desirable to release older prisoners from jail. Second, as a matter of logic, the incapacitative rationale might imply that a person should be put in jail even if he has not committed a crime — if his danger to society makes incapacitating him

⁷⁸ For simplicity, we are not taking into account here the gains that individuals obtain from committing offenses or the disutility that individuals bear from time in jail.

⁷⁹ See, for example, Wilson and Herrnstein (1985, pp. 126-147) and U.S. Department of Justice (1997a, pp. 371, table 4.4, 378-379, table 4.7).

worthwhile. This would be true, for example, if there were some accurate way to predict a person's dangerousness independently of his actual behavior. In practice, however, the fact that a person has committed a harmful act may be a good basis for predicting his future behavior, in which case the incapacitation rationale would imply that a jail term should be imposed only if the individual has committed an especially harmful act.

The optimal probability of detection is determined by a straightforward tradeoff. The higher the probability, the greater the number of individuals who will be incapacitated, resulting in social gains equal to the difference between the harm that individuals would cause and the cost of their incapacitation. But the higher the probability, the higher are enforcement costs. At some point, it is optimal to stop raising the probability, when the marginal social gains just equal the marginal cost of raising the probability.

Last, we briefly comment on the relationship between the nature of optimal enforcement when incapacitation is the goal versus when deterrence is the goal. First, when incapacitation is the goal, the optimal length of the prison term (which is determined by the condition that $h(t) > c$) is independent of the probability of apprehension. In contrast, when deterrence is the goal, the optimal sanction depends on the probability — the sanction generally is higher the lower is the probability. Second, when incapacitation is the goal, the probability and magnitude of sanctions are independent of the ability to deter. Thus, for example, if this ability is limited (consider individuals who commit crimes while enraged), a low expected sanction may be optimal under the deterrence rationale, but a high expected sanction still might be called for to incapacitate.⁸⁰

⁸⁰ See Shavell (1987a) for a theoretical examination of optimal incapacitation policy, Ehrlich (1981, pp. 315-316, 319-321) for a model used to estimate the relative importance of incapacitation and deterrence, and Levitt (1998) and Kessler and Levitt (1999) for empirical studies of incapacitation and deterrence. Economists have paid much less attention to incapacitation than to deterrence, despite the significance of the incapacitation rationale in criminal law enforcement.

26. COSTLY OBSERVATION OF WEALTH

In our prior discussions of optimal sanctions, we implicitly assumed that the enforcement authority could costlessly observe individuals' wealth levels. Knowing wealth levels, the enforcement authority then chose the sanctions to impose, fines and/or imprisonment sentences. In fact, however, individuals and firms may be able to hide assets from government enforcers, including by hoarding cash, transferring assets to relatives or related legal entities, or moving money to offshore bank accounts. In this section we consider optimal sanctions when an individual's level of wealth can be observed only after a costly audit or not at all.⁸¹

Suppose first that the enforcement authority employs fines as sanctions and can audit an individual who claims that he cannot pay the fine. If the audit determines that the individual misrepresented his wealth level, he can be fined for having lied about his wealth. Assume for simplicity that this fine is independent of the degree of misrepresentation. The enforcement authority's problem is to choose, so as to maximize social welfare, the probability of detecting the offense, the fine for the offense, the probability of an audit conditional on the individual's claiming that he cannot pay the fine, and the fine for misrepresentation of wealth. Without loss of generality, we assume that if the latter fine is applicable, it is imposed instead of the fine for the offense (rather than in addition to the fine for the offense).

It can be demonstrated that the optimal fine for misrepresenting one's wealth level equals the fine for the offense divided by the audit probability, and therefore generally exceeds the fine for the offense. This is a natural generalization of the formula for the optimal fine, given the probability of detection, which is the harm divided by the probability. In effect, the "harm" from

⁸¹ The discussion in this section is based on Polinsky (2004a; 2004b). See also Chu and Jiang (1993) and Levitt (1997), who consider the choice between fines and imprisonment when wealth cannot be discovered by the enforcement authority at any cost (in Chu and Jiang's case, this assumption is implicit), and Garoupa (1998), who investigates optimal fines when the enforcement authority is assumed to costlessly observe an underestimate of offenders' wealth levels.

the act of concealing one's wealth is the failure to pay the fine for the original offense, so the optimal fine for concealment is this harm divided by the applicable probability of being caught if one engages in it.

Assuming the harmful act is worth controlling, the optimal audit probability is positive, increases as the cost of an audit declines, and equals one if the cost is sufficiently low. Auditing is valuable because it reduces misrepresentation of wealth and thereby increases deterrence. If the optimal audit probability is less than one, however, some individuals who are capable of paying the fine for the offense will misrepresent their wealth levels. Unlike in the basic analysis in which wealth is assumed to be costlessly observable, the optimal fine for the offense now results in underdeterrence, due to the cost of auditing wealth levels. By reducing the fine for the offense from the level that would lead to first-best behavior, fewer individuals will misrepresent their wealth levels, so auditing costs decline. The reduction in deterrence has no first-order effect on social welfare because the marginal individuals who are induced to commit the offense have gains equal to the harm from the offense.

Next suppose that the enforcement authority simply cannot observe wealth, say because the cost of an audit is prohibitively high.⁸² If the enforcement authority would have used fines alone if it could have observed wealth at no cost, it would have imposed a higher fine on higher-wealth individuals.⁸³ It obviously cannot do this now. Instead, it may be desirable to use the threat of an imprisonment sentence to induce individuals capable of paying a higher fine to do so. Specifically, if there are two levels of wealth, the enforcement authority might set the fine

⁸² The following discussion introduces imprisonment as an alternative, or supplement, to fines. For simplicity, imprisonment was not considered above in the discussion of auditing.

⁸³ For now familiar reasons, it would have done this to reduce or eliminate the underdeterrence that otherwise would occur if the fine were the same as for low-wealth individuals.

greater than the wealth level of the low-wealth individuals and impose an imprisonment sentence on any individuals who do not pay this fine. By using an imprisonment sentence in this way, high-wealth individuals can be induced to pay a higher fine, which is socially beneficial, though low-wealth individuals now will incur a socially costly sanction, imprisonment. In other words, when wealth is not observable, it may be desirable to impose a costly sanction — imprisonment sentences — on low-wealth individuals in order to better deter high-wealth offenders through a cheap sanction — fines.

Another possibility is that the enforcement authority would have used both fines and imprisonment if it could have observed wealth at no cost. Perhaps surprisingly, the inability to observe wealth is of no consequence in this case. The reason, in essence, is that the mix of fines and imprisonment that would be chosen when wealth is observable will impose a higher burden (though a lower fine) on low-wealth individuals. Thus, high-wealth individuals will naturally want to identify themselves. Specifically, they will prefer to pay a higher fine and bear a shorter imprisonment sentence than to masquerade as low-wealth individuals, who will bear longer imprisonment sentences and a higher overall burden. Consequently, the same mix of sanctions that would have been imposed on both groups if wealth were costlessly observable can be used when wealth is not observable.

In summary, information about wealth levels is useful only if the enforcement authority would want to impose a higher burden of sanctions on high-wealth individuals than on low-wealth individuals, for then high-wealth individuals would pretend to be low-wealth individuals if wealth could not be observed. This is the case if the enforcement authority would want to use fines alone. If, however, the enforcement authority would want to impose a lower burden of sanctions on high-wealth individuals than on low-wealth individuals, high-wealth individuals will voluntarily bear such sanctions even if they include a higher fine. This case is applicable

when the enforcement authority would want to use fines and imprisonment together. Monitoring of wealth levels may be worthwhile in the first case, but is not needed in the second case.

27. SOCIAL NORMS

Although we have restricted attention to public enforcement of law, social norms and morality should be mentioned because they influence in significant ways the attainment of desired behavior.⁸⁴ By social norms (or moral rules), we mean conduct — such as keeping promises, not lying, and not harming others — that is associated with certain distinctive psychological and social attributes. In particular, social norms influence behavior partly through internal incentives: when a person obeys a moral rule, he will tend to feel virtuous, and if he disobeys the rule, he will tend to feel guilty. Social norms also influence behavior through external incentives: when a person is observed by another party to have obeyed a moral rule, that party may bestow praise on the first party, who will enjoy the praise; and if the person is observed by the other party to have disobeyed the rule, the second party will tend to disapprove of the first party, who will dislike the disapproval.

Because social norms affect behavior through the foregoing moral incentives, some socially desirable conduct can be encouraged reasonably well without employing the legal system.⁸⁵ For example, whether an individual cuts in line at the movie theater, keeps his lunch engagements, or lets his children make a nuisance of themselves at the supermarket, is controlled with rough success by internal and external moral incentives. Such conduct generally can be regulated satisfactorily by moral incentives alone because, among other things, internal

⁸⁴ On social norms and the law, see generally McAdams and Rasmusen (2006). See also University of Pennsylvania Law Review (1996) and Posner (1997).

⁸⁵ For a comparison of social norms and law enforcement as means of controlling behavior, see Shavell (2002).

incentives work automatically (we know when we have done the wrong thing), external incentives are likely to apply (if a person cuts in line, this will be noticed), and the benefits from violations are not large, so that they can be outweighed relatively easily by the force of the moral incentives.

The need for formal law enforcement stems from two principal considerations. First, much conduct that society desires cannot be controlled through moral incentives alone. One reason is that the private gains from undesirable conduct are often large. The utility obtained by a robber, a tax cheat, or a polluter may be substantial, and dominate the moral incentives. Another reason is that external moral sanctions might be imposed only with a low probability (the robber, tax cheat, or polluter might not be spotted by others). A second rationale for formal law enforcement is that the social harm from failing to control an act through moral incentives may be large. This makes the expense of law enforcement worth incurring (as in the case of robbery, but not of cutting in line at movie theaters).

Although we have been treating social norms and formal law enforcement as distinct ways of controlling behavior, law enforcement might also influence social norms.⁸⁶ For instance, enforcement of laws prohibiting discrimination based on race may change beliefs about proper conduct (reinforcing internal moral incentives) and lead to a greater willingness of individuals to express disapproval when they witness discriminatory behavior (enhancing external moral incentives). The importance of the effect of law on social norms may be limited, however, to the extent that social norms are mainly the result of early childhood experience and the messages conveyed by parents and other authority figures, such as educators.

⁸⁶ See, for example, McAdams (1997) and Sunstein (1996).

28. FAIRNESS

To this point we have not considered the possibility that individuals have opinions about the fairness of sanctions or the arbitrariness of enforcement.⁸⁷

Suppose, first, that individuals believe that sanctions should be imposed on those who have committed certain bad acts and that the magnitude of sanctions should reflect the gravity of the acts. A formal assumption that captures this view is that individuals obtain fairness-related utility from the imposition of sanctions on those who committed the bad acts, where this utility is maximized at a *fairness-ideal* level of sanction that depends on the harmfulness (or a related aspect) of the acts. The fairness-ideal sanction may be lower or higher than the conventionally optimal sanction that we have discussed above. Note in particular that the fairness-ideal sanction does not depend on the probability of detection, whereas the conventionally optimal sanction is higher the lower is the probability of detection, suggesting that if the probability of detection is sufficiently low, the conventionally optimal sanction will exceed the fairness-ideal sanction. In any case, when fairness-related utility is taken into account along with the other elements of social welfare considered above, the optimal sanction will be a compromise between the fairness-ideal sanction and the conventionally optimal sanction.

When both the probability and magnitude of sanctions may be varied, the conventional solution to the enforcement problem also is altered because of fairness considerations. As discussed previously, if individuals are risk neutral, the usual solution consists of the highest possible sanction and a relatively low probability. When the issue of fairness is added to the

⁸⁷ The discussion in this section is based on Polinsky and Shavell (2000b) and Kaplow and Shavell (2002, pp. 291-378), the latter of which relates the economic analysis of fairness in enforcement to the philosophical literature. See also Miceli (1991) and Diamond (2002), who derive optimal sanctions taking both deterrence and the fairness of sanctions into account, but holding the probability of sanctions fixed. Waldfogel (1993) studies empirically whether actual sanctions are better explained by considerations of deterrence or fairness.

analysis, however, the usual solution generally is not optimal because a very high sanction will be seen as unfair, or more precisely, will result in the lowering of individuals' fairness-related utility. With respect to double parking, for example, even a sanction of \$100 might be considered unfair because double parking is regarded as only a slightly bad act.

A consequence of the desire to keep sanctions at fair levels, meaning quite constrained levels for acts that are not very bad or harmful, is that the socially optimal probability of detection changes. One possibility is that the optimal probability would be higher, perhaps much higher than the conventionally optimal probability: to achieve a desired level of deterrence with a lower fairness-restricted sanction, the probability has to rise. If the sanction for double parking cannot exceed \$100 because of fairness considerations, then, to create an expected sanction of \$10, the probability must be 10%, greatly exceeding the approximately 1/10% probability that would be optimal if risk-neutral individuals have wealth of \$10,000 and the fine is set at this level. Another possibility, though, is that the optimal probability would be lower than in the conventional case: the additional deterrence created by raising the probability might be relatively low because the sanction is relatively low; and the lower the deterrent benefit from raising the probability, the lower would be the social incentive to devote resources to enforcement.

Another concept of fairness concerns the probability of detection rather than the magnitude of sanctions. Suppose that individuals consider it unfair for some violators of law to be sanctioned when others who were lucky enough not to be caught are not sanctioned. Specifically, suppose that individuals experience fairness-related disutility if there is only a probability rather than a certainty of sanctions, and their disutility rises the lower the probability of detection. Then the optimal probability would be higher, and therefore the optimal sanction would be lower, than in the absence of this fairness-related component of social welfare.

A further notion of fairness involves the form of liability, whether liability is strict or

based on fault. Individuals might prefer fault-based liability because sanctions are imposed on parties only if they behaved in a socially inappropriate way. If individuals derive greater fairness-related utility from use of fault-based liability, then this form of liability is more likely to be optimal than we have suggested previously.

A final issue concerns the relevance of fairness considerations when firms, as opposed to individuals, are sanctioned. If sanctions are imposed on firms, then fairness-related utility may have to be reconsidered, presuming that what matters in terms of fairness is that the *individuals* responsible for harmful acts bear sanctions as opposed to the artificial legal entity of a firm. Specifically, one would want to identify the sanctions actually suffered by such persons within a firm if the firm bears a sanction. For example, if a firm demotes a person who negligently caused the firm to incur a sanction, then the person's loss from the demotion would be the fairness-relevant sanction, not the sanction imposed on the firm. Note, too, that the imposition of sanctions on firms often penalizes individuals who are unlikely to be considered responsible for the harm, namely shareholders and customers. To the extent that the fairness goal is to penalize responsible individuals within firms, and not firms as entities, the social value of sanctions in achieving fairness is attenuated.

29. CONCLUSION

Having reviewed the theory of public enforcement of law, we want to conclude by commenting on two types of private behavior that bear significantly on public law enforcement.

First, private parties may themselves take actions to prevent being harmed. For example, to reduce the risk of being criminally victimized, individuals might install locks on their possessions, carry weapons, or hire security personnel. These private efforts to prevent or deter harmful acts serve as a partial substitute for public efforts; moreover, private efforts are

sometimes more efficient than public efforts (citizens may know better where to put locks), though they also may be less efficient (public authorities may know better how to assign police). These observations raise important questions about the relationship between private protection and public enforcement. Should private efforts to protect against being victimized be regulated? Does the state spend too little on public enforcement, relying on the fact that private actors often undertake their own defensive efforts? The optimal relationship between private and public efforts to control harmful activities deserves more careful examination.⁸⁸

Second, private individuals may bring suits against parties who also may be subject to publicly imposed penalties. For example, a victim of an automobile accident may sue the driver who caused the accident, and this driver might also be sanctioned by the state for a driving infraction. Private lawsuits channel harm-creating behavior and thus constitute a substitute, at least to some extent, for public enforcement. This leads one to ask how public enforcement and parallel private litigation should be managed. Should the payment of a public penalty be an offset to private damages, and vice versa? Should the state regulate private litigation so as to better coordinate it with public enforcement? Is public enforcement or private litigation the socially cheaper way to accomplish desired behavior?

A full treatment of the control of harm-creating behavior would address both sets of issues just discussed.

⁸⁸ There is some literature that discusses the issues raised in this paragraph; see, for example, Clotfelter (1977; 1978) and Shavell (1991a).

REFERENCES

- Andenaes, J. (1966), "The General Preventive Effects of Punishment", *University of Pennsylvania Law Review* 114: 949-983.
- Andreoni, J. (1991), "Reasonable Doubt and the Optimal Magnitude of Fines: Should the Penalty Fit the Crime?" *RAND Journal of Economics* 22: 385-395.
- Andreoni, J., B. Erard and J. Feinstein (1998), "Tax Compliance", *Journal of Economic Literature* 36: 818-860.
- Arlen, J.A. (1994), "The Potentially Perverse Effects of Corporate Criminal Liability", *Journal of Legal Studies* 23: 833-867.
- Bardhan, P. (1997), "Corruption and Development: A Review of Issues", *Journal of Economic Literature* 35: 1320-1346.
- Bar-Gill, O. and A. Harel (2001), "Crime Rates and Expected Sanctions: The Economics of Deterrence Revisited", *Journal of Legal Studies* 30: 485-501.
- Bebchuk, L.A. and L. Kaplow (1992), "Optimal Sanctions When Individuals Are Imperfectly Informed About the Probability of Apprehension", *Journal of Legal Studies* 21: 365-370.
- Bebchuk, L.A. and L. Kaplow (1993), "Optimal Sanctions and Differences in Individuals' Likelihood of Avoiding Detection", *International Review of Law and Economics* 13: 217-224.
- Beccaria, C. (1767), *On Crimes and Punishments, and Other Writings* (Cambridge University Press, New York, Editor, R. Bellamy, Translator, R. Davies et al., 1995).
- Becker, G.S. (1968), "Crime and Punishment: An Economic Approach", *Journal of Political Economy* 76: 169-217.
- Becker, G.S. and G.J. Stigler (1974). "Law Enforcement, Malfeasance, and Compensation of Enforcers", *Journal of Legal Studies* 3: 1-18.

- Bentham, J. (1789), *An Introduction to the Principles of Morals and Legislation*, in: *The Utilitarians* (Anchor Books, Garden City, N.Y., 1973).
- Block, M.K. and J.G. Sidak (1980), “The Cost of Antitrust Deterrence: Why Not Hang a Price Fixer Now and Then?”, *Georgetown Law Journal* 68: 1131-1139.
- Bouckaert, B. and G. De Geest, eds. (1992), *Bibliography of Law and Economics* (Kluwer Academic Publishers, Dordrecht).
- Bowles, R. and N. Garoupa (1997), “Casual Police Corruption and the Economics of Crime”, *International Review of Law and Economics* 17: 75-87.
- Brown, J.P. (1973), “Toward an Economic Theory of Liability”, *Journal of Legal Studies* 2: 323-349.
- Carr-Hill, R.A. and N. H. Stern (1979), *Crime, the Police and Criminal Statistics* (Academic Press, London).
- Chu, C.Y.C., S. Hu and T. Huang (2000), “Punishing Repeat Offenders More Severely”, *International Review of Law and Economics* 20: 127-140.
- Chu, C.Y.C. and N. Jiang (1993), “Are Fines More Efficient than Imprisonment?”, *Journal of Public Economics* 51: 391-413.
- Clotfelter, C.T. (1977), “Public Services, Private Substitutes, and the Demand for Protection Against Crime”, *American Economic Review* 67: 867-877.
- Clotfelter, C.T. (1978), “Private Security and the Public Safety”, *Journal of Urban Economics* 5: 388-402.
- Cohen, M.A. (1989), “Corporate Crime and Punishment: A Study of Social Harm and Sentencing Practice in the Federal Courts, 1984-1987”, *American Criminal Law Review* 26: 605-660.
- Cohen, M.A. (1999), “Monitoring and Enforcement of Environmental Policy”, in: H. Folmer and

- T. Tietenberg, eds., *The International Yearbook of Environmental and Resource Economics 1999/2000* (Edward Elgar, Cheltenham, United Kingdom): 44-106.
- Cooter, R.D. and D.L. Rubinfeld (1989), "Economic Analysis of Legal Disputes and Their Resolution", *Journal of Economic Literature* 27: 1067-1097.
- Craswell, R. and J.E. Calfee (1986), "Deterrence and Uncertain Legal Standards", *Journal of Law, Economics, & Organization* 2: 279-303.
- Diamond, P.A. (2002), "Integrating Punishment and Efficiency Concerns in Punitive Damages for Reckless Disregard of Risks to Others", *Journal of Law, Economics, & Organization* 18: 117-139.
- Eide, E. (2000), "Economics of Criminal Behavior", in: B. Bouckaert and G. De Geest, eds., *Encyclopedia of Law and Economics* (Edward Elgar Publishing Limited, Cheltenham, United Kingdom): 345-389.
- Ehrlich, I. (1981), "On the Usefulness of Controlling Individuals: An Economic Analysis of Rehabilitation, Incapacitation and Deterrence", *American Economic Review* 71: 307-332.
- Emons, W. (2003), "A Note on the Optimal Punishment for Repeat Offenders", *International Review of Law and Economics* 23: 253-259.
- Friedman, D.D. (1981), "Reflections on Optimal Punishment, or: Should the Rich Pay Higher Fines?" *Research in Law and Economics* 3: 185-205.
- Friedman, D.D. (1995), "Making Sense of English Law Enforcement in the Eighteenth Century", *University of Chicago Law School Roundtable* 2: 475-505.
- Garoupa, N. (1997), "The Theory of Optimal Law Enforcement", *Journal of Economic Surveys* 11: 267-295.
- Garoupa, N. (1998), "Optimal Law Enforcement and Imperfect Information when Wealth Varies Among Individuals", *Economica* 65: 479-490.

- Garoupa, N. and D. Klerman (2004), "Corruption and the Optimal Use of Nonmonetary Sanctions", *International Review of Law and Economics* 24: 219-225.
- Glaeser, E.L., D.P. Kessler and A.M. Piehl (2000), "What Do Prosecutors Maximize? An Analysis of the Federalization of Drug Crimes", *American Law and Economics Review* 2: 259-290.
- Grossman, G.M. and M.L. Katz (1983), "Plea Bargaining and Social Welfare", *American Economic Review* 73: 749-757.
- Hindriks, J., M. Keen and A. Muthoo (1999), "Corruption, Extortion and Evasion", *Journal of Public Economics* 74: 395-430.
- Innes, R. (1999), "Remediation and Self-Reporting in Optimal Law Enforcement", *Journal of Public Economics* 72: 379-393.
- Kahneman, D., P. Slovic and A. Tversky, eds. (1982), *Judgment under Uncertainty: Heuristics and Biases* (Cambridge University Press, Cambridge and New York).
- Kaplow, L. (1990a), "Optimal Deterrence, Uninformed Individuals, and Acquiring Information About Whether Acts are Subject to Sanctions", *Journal of Law, Economics, & Organization* 6: 93-128.
- Kaplow, L. (1990b), "A Note on the Optimal Use of Nonmonetary Sanctions", *Journal of Public Economics* 42: 245-247.
- Kaplow, L. (1992), "The Optimal Probability and Magnitude of Fines for Acts That Definitely Are Undesirable", *International Review of Law and Economics* 12: 3-11.
- Kaplow, L. and S. Shavell (1994a), "Accuracy in the Determination of Liability", *Journal of Law and Economics* 37: 1-15.
- Kaplow, L. and S. Shavell (1994b), "Optimal Law Enforcement with Self-Reporting of Behavior", *Journal of Political Economy* 102: 583-606.

- Kaplow, L. and S. Shavell (2002), *Fairness Versus Welfare* (Harvard University Press, Cambridge).
- Kenkel, D.S. (1993), “Do Drunk Drivers Pay Their Way? A Note on Optimal Penalties for Drunk Driving”, *Journal of Health Economics* 12: 137-149.
- Kessler, D. and S.D. Levitt (1999), “Using Sentence Enhancements to Distinguish between Deterrence and Incapacitation”, *Journal of Law and Economics* 42: 343-363.
- Klitgaard, R. (1988), *Controlling Corruption* (University of California Press, Berkeley).
- Kobayashi, B.H. and J.R. Lott, Jr. (1992), “Low-Probability-High-Penalty Enforcement Strategies and the Efficient Operation of the Plea-Bargaining System”, *International Review of Law and Economics* 12: 69-77.
- Kornhauser, L.A. (1982), “An Economic Analysis of the Choice Between Enterprise and Personal Liability for Accidents”, *California Law Review* 70: 1345-1392.
- Laffont, J. and J. Tirole (1993), *The Theory of Incentives in Procurement and Regulation* (MIT Press, Cambridge, MA).
- Landes, W.M. (1971), “An Economic Analysis of the Courts”, *Journal of Law and Economics* 4: 61-107.
- Landes, W.M. and R.A. Posner (1975), “The Private Enforcement of Law”, *Journal of Legal Studies* 4: 1-46.
- 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. (1998), “Why Do Increased Arrest Rates Appear to Reduce Crime: Deterrence, Incapacitation, or Measurement Error?” *Economic Inquiry* 36: 353-372.
- Levitt, S.D. and T. Miles (2006), “Empirical Study of Criminal Punishment”, in: A.M. Polinsky and S. Shavell, eds., *Handbook of Law and Economics*, vol. 1 (North-Holland,

- Amsterdam).
- Lewin, J.L. and W.N. Trumbull (1990), "The Social Value of Crime?" *International Review of Law and Economics* 10: 271-284.
- McAdams, R.H. (1997), "The Origin, Development, and Regulation of Norms", *Michigan Law Review* 96: 338-433.
- McAdams, R.H., and E. Rasmusen (2006), "Norms in Law and Economics", in: A.M. Polinsky and S. Shavell, eds., *Handbook of Law and Economics*, vol. 2 (North-Holland, Amsterdam).
- Malik, A.S. (1990), "Avoidance, Screening and Optimum Enforcement", *RAND Journal of Economics* 21: 341-353.
- Malik, A.S. (1993), "Self-Reporting and the Design of Policies for Regulating Stochastic Pollution", *Journal of Environmental Economics and Management* 24: 241-257.
- Miceli, T.J. (1991), "Optimal Criminal Procedure: Fairness and Deterrence", *International Review of Law and Economics* 11: 3-10.
- Miceli, T.J. (1996), "Plea Bargaining and Deterrence: An Institutional Approach", *European Journal of Law and Economics* 3: 249-264.
- Miceli, T.J. and C. Bucci (2005), "A Simple Theory of Increasing Penalties for Repeat Offenders", *Review of Law & Economics* 1: 71-80 (available at <http://www.bepress.com/rle/vol1/iss1/art5>).
- Montesquieu (1748), *The Spirit of the Laws* (University of California Press, Berkeley, Repr. ed. 1977).
- Mookherjee, D. (1997), "The Economics of Enforcement", in: A. Bose, M. Rakshit and A. Sinha, eds., *Issues in Economic Theory and Policy, Essays in Honor of Tapas Majumdar* (Oxford University Press, New Delhi) 202-249.

- Mookherjee, D. and I.P.L. Png (1992), "Monitoring vis-à-vis Investigation in Enforcement of Law", *American Economic Review* 82: 556-565.
- Mookherjee, D. and I.P.L. Png (1994), "Marginal Deterrence in Enforcement of Law", *Journal of Political Economy* 102: 1039-1066.
- Mookherjee, D. and I.P.L. Png (1995), "Corruptible Law Enforcers: How Should They Be Compensated?" *Economic Journal* 105: 145-159.
- Newman, H.A. and D.W. Wright (1990), "Strict Liability in a Principal-Agent Model", *International Review of Law and Economics* 10: 219-231.
- Pashigian, B.P. (1975), "On the Control of Crime and Bribery", *Journal of Legal Studies* 4: 311-326.
- Png, I.P.L. (1986), "Optimal Subsidies and Damages in the Presence of Judicial Error", *International Review of Law and Economics* 6: 101-105.
- Polinsky, A.M. (1980a), "Private versus Public Enforcement of Fines", *Journal of Legal Studies* 9: 105-127.
- Polinsky, A.M. (1980b), "Strict Liability vs. Negligence in a Market Setting", *American Economic Review: Papers and Proceedings* 70: 363-370.
- Polinsky, A.M. (2004a), "Optimal Fines and Auditing When Wealth is Costly to Observe", Working Paper No. 289, John M. Olin Program in Law and Economics, Stanford Law School.
- Polinsky, A.M. (2004b), "The Optimal Use of Fines and Imprisonment When Wealth is Unobservable", Working Paper No. 290, John M. Olin Program in Law and Economics, Stanford Law School (forthcoming in the *Journal of Public Economics*).
- Polinsky, A.M. and D.L. Rubinfeld (1988), "The Deterrent Effects of Settlements and Trials", *International Review of Law and Economics* 8: 109-116.

- Polinsky, A.M. and D.L. Rubinfeld (1989), "A Note on Optimal Public Enforcement with Settlements and Litigation Costs", *Research in Law and Economics* 12: 1-8.
- Polinsky, A.M. and D.L. Rubinfeld (1991), "A Model of Optimal Fines for Repeat Offenders", *Journal of Public Economics* 46: 291-306.
- 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 (1991), "A Note on Optimal Fines When Wealth Varies Among Individuals", *American Economic Review* 81: 618-621.
- Polinsky, A.M. and S. Shavell (1992), "Enforcement Costs and the Optimal Magnitude and Probability of Fines", *Journal of Law and Economics* 35: 133-148.
- Polinsky, A.M. and S. Shavell (1993), "Should Employees Be Subject to Fines and Imprisonment Given the Existence of Corporate Liability?" *International Review of Law and Economics* 13: 239-257.
- Polinsky, A.M. and S. Shavell (1998), "On Offense History and the Theory of Deterrence", *International Review of Law and Economics* 18: 305-324.
- Polinsky, A.M. and S. Shavell (1999), "On the Disutility and Discounting of Imprisonment and the Theory of Deterrence", *Journal of Legal Studies* 28: 1-16.
- Polinsky, A.M. and S. Shavell (2000a), "The Economic Theory of Public Enforcement of Law", *Journal of Economic Literature* 38: 45-76.
- Polinsky, A.M. and S. Shavell (2000b), "The Fairness of Sanctions: Some Implications for Optimal Enforcement Policy", *American Law and Economics Review* 2: 223-237.
- Polinsky, A.M. and S. Shavell (2001), "Corruption and Optimal Law Enforcement", *Journal of*

- Public Economics* 81: 1-24.
- Posner, R.A. (1980a), "Optimal Sentences for White-Collar Criminals", *American Criminal Law Review* 17: 409-418.
- Posner, R.A. (1980b), "Retribution and Related Concepts of Punishment", *Journal of Legal Studies* 9: 71-92.
- Posner, R.A. (1985), "An Economic Theory of the Criminal Law", *Columbia Law Review* 85: 1193-1231.
- Posner, R.A. (1997), "Social Norms and the Law: An Economic Approach", *American Economic Review: Papers and Proceedings* 87: 365-369.
- Reinganum, J.F. (1988), "Plea Bargaining and Prosecutorial Discretion", *American Economic Review* 78: 713-728.
- Rose-Ackerman, S. (1999), *Corruption and Government: Causes, Consequences, and Reform* (Cambridge University Press, Cambridge).
- Rubinstein, A. (1979), "An Optimal Conviction Policy for Offenses that May Have Been Committed by Accident", in: S. J. Brams, A. Schotter, and G. Schwödiauer, eds., *Applied Game Theory* (Physica-Verlag, Wurzburg) 406-413.
- Sah, R.K. (1991), "Social Osmosis and Patterns of Crime", *Journal of Political Economy* 99: 1272-1295.
- Segerson, K. and T. Tietenberg (1992), "The Structure of Penalties in Environmental Enforcement: An Economic Analysis", *Journal of Environmental Economics and Management* 23: 179-200.
- Shavell, S. (1980), "Strict Liability versus Negligence", *Journal of Legal Studies* 9: 1-25.
- Shavell, S. (1982), "On Liability and Insurance", *Bell Journal of Economics* 13: 120-132.
- Shavell, S. (1985), "Criminal Law and the Optimal Use of Nonmonetary Sanctions as a

- Deterrent”, *Columbia Law Review* 85: 1232-1262.
- Shavell, S. (1987a), “A Model of Optimal Incapacitation”, *American Economic Review: Papers and Proceedings* 77: 107-110.
- Shavell, S. (1987b), “The Optimal Use of Nonmonetary Sanctions as a Deterrent”, *American Economic Review* 77: 584-592.
- Shavell, S. (1991a), “Individual Precautions to Prevent Theft: Private versus Socially Optimal Behavior”, *International Review of Law and Economics* 11: 123-132.
- Shavell, S. (1991b), “Specific versus General Enforcement of Law”, *Journal of Political Economy* 99: 1088-1108.
- Shavell, S. (1992), “A Note on Marginal Deterrence”, *International Review of Law and Economics* 12: 345-355.
- Shavell, S. (1993), “The Optimal Structure of Law Enforcement”, *Journal of Law and Economics* 36: 255-287.
- Shavell, S. (1997a), “The Fundamental Divergence Between the Private and the Social Motive to Use the Legal System”, *Journal of Legal Studies* 26: 575-612.
- Shavell, S. (1997b), “The Optimal Level of Corporate Liability Given the Limited Ability of Corporations to Penalize Their Employees”, *International Review of Law and Economics* 17: 203-213.
- Shavell, S. (2002), “Law versus Morality as Regulators of Conduct”, *American Law and Economics Review* 4: 227-257.
- Shavell, S. (2004), *Foundations of Economic Analysis of Law* (Harvard University Press, Cambridge, MA).
- Shleifer, A. and R.W. Vishny (1993), “Corruption”, *Quarterly Journal of Economics* 108: 599-617.

- Spier, K.E. (1997), “A Note on the Divergence Between the Private and the Social Motive to Settle Under a Negligence Rule”, *Journal of Legal Studies* 26: 613-621.
- Spier, K.E. (2006), “Litigation”, in: A.M. Polinsky and S. Shavell, eds., *Handbook of Law and Economics*, vol. 1 (North-Holland, Amsterdam).
- Stigler, G.J. (1970), “The Optimum Enforcement of Laws”, *Journal of Political Economy* 78: 526-536.
- Sunstein, C.R. (1996), “Social Norms and Social Roles”, *Columbia Law Review* 96: 201-266.
- Sykes, A.O. (1981), “An Efficiency Analysis of Vicarious Liability Under the Law of Agency”, *Yale Law Journal* 91: 168-206.
- University of Pennsylvania Law Review (1996), “Symposium: Law, Economics, and Norms”, *University of Pennsylvania Law Review* 144: 1643-2339.
- U.S. Department of Justice (1988), *Report to the Nation on Crime and Justice, Technical Appendix* (U.S. Department of Justice, Bureau of Justice Statistics, 2nd ed., NCJ-112011, Washington, D.C.).
- U.S. Department of Justice (1997a), *Sourcebook of Criminal Justice Statistics — 1996* (U.S. Department of Justice, Office of Justice Statistics, Bureau of Justice Statistics, NCJ-165361, Washington, D.C.).
- U.S. Department of Justice (1997b), *Uniform Crime Reports for the United States, 1996* (U.S. Department of Justice, Federal Bureau of Investigation, Washington, D.C.).
- U.S. Department of Justice (1998), *Profile of Jail Inmates 1996* (U.S. Department of Justice, Office of Justice Programs, NCJ-164620, Washington D.C.).
- Waldfogel, J. (1993), “Criminal Sentences as Endogenous Taxes: Are They ‘Just’ or ‘Efficient’?”, *Journal of Law and Economics* 36: 139-151.
- Waldfogel, J. (1993), “Criminal Sentences as Endogenous Taxes: Are They “Just” or

“Efficient”?”, *Journal of Law and Economics* 36: 139-151.

Wilde, L.L. (1992), “Criminal Choice, Nonmonetary Sanctions, and Marginal Deterrence: A Normative Analysis”, *International Review of Law and Economics* 12: 333-344.

Wilson, J.Q. and R.J. Herrnstein (1985), *Crime and Human Nature* (Simon and Schuster, New York).

Zimring, F.E. and G.J. Hawkins (1973), *Deterrence* (University of Chicago Press, Chicago).