# THE UNIVERSITY of EDINBURGH

This thesis has been submitted in fulfilment of the requirements for a postgraduate degree (e.g. PhD, MPhil, DClinPsychol) at the University of Edinburgh. Please note the following terms and conditions of use:

# Essays on Corruption and Inequality

*Christopher Stapenhurst*



*Doctor of Philosophy*

THE UNIVERSITY OF EDINBURGH

2022

*To my departed friend and mentor,*

*Peter McCaffery.*

# Abstract

The first part of my thesis studies a problem in which a polluting firm can bribe an inspector to conceal evidence of illegal pollution. The government can deter bribes by rewarding the inspector to report the evidence, but doing so is costly. How can we reduce these costs?

Chapter 1, "Two Corruptible Monitors", asks whether more inspectors help to deter bribes. I compare optimal solutions for moral hazard problems in which the government pays two inspectors to provide hard evidence about the firm's emissions. If the firm cannot bribe both monitors simultaneously, then the government might prefer to have each inspector receive a different piece of evidence, because then the firm can never conceal all of it. But if the firm can bribe both of the inspectors, then the government always prefers to have one inspector receive all the evidence. Why? Because when each inspector receives a different piece of evidence, the government must pay both inspectors in order to deter bribes to conceal one or other piece. But when one inspector receives both pieces of evidence, the same payment that deters bribes to conceal one piece, also deters bribes to conceal the other piece, so the government pays less overall. I conclude that coalition formation frictions play an important role in incentive design problems with more than one inspector.

Chapter 2, "Lemons by design: sowing secrets that curb corruption", with Andrew Clausen, studies one particular coalition formation friction: asymmetric information. We find that the regulator can artificially create distrust between the firm and the inspector by promising to pay a secret reward to one or the other whenever the inspector reports incriminating evidence; and by giving them each different clues about who stands to be rewarded. We find that the cheapest way to deter bribes is (i) to secretly select either the firm or the inspector to 'win' a reward whenever evidence is reported; and (ii) to give both the firm and the inspector a secret clue about who will win. If the inspector conceals the evidence, then the winner forgoes their reward — i.e. they 'get a lemon'. The distribution of clues is carefully constructed to engineer the worst possible lemons problem in the market for concealment: player $i$ only enters the market if her clue is strong enough to make her believe that player $j$ is the winner, despite knowing that player $j$ only enters if his clue indicates that player $i$ is the winner. But then higher order reasoning leads neither player to enter the market, no matter what clue they receive. Hence, bribery never takes place in equilibrium. As well as deterring bribes cheaply and robustly, this result demonstrates the full extent of contagious adverse selection in bilateral trades. In particular, it proves that there exist two-sided adverse selection problems which are strictly worse (in the sense of destroying surplus) than the worse possible one-sided adverse selection problem.

The second part of my thesis studies the measurement of inequality in ordinal variables such as self-reported wellbeing, sanitation, and perceived corruption. The median-preserving spread (MPS) partial ordering for ordinal variables (Alison and Foster 2004) has become ubiquitous in inequality measurement, but the literature lacks a frequentist method for inferring that one population is more unequal than another according to the MPS criterion. Chapter 3, "Inferring Inequality: Testing for Median-Preserving Spreads in Ordinal Data" with Ramses Abul Naga and Gaston Yalonetzky, devises statistical tests of the hypothesis that a distribution $G$ is not an MPS of a distribution $F$. Rejecting this hypothesis leads to the conclusion that $G$ is more unequal than $F$ according to the MPS criterion. Monte Carlo simulations and novel graphical techniques show that a simple, asymptotic Z test is sufficient for most applications. We illustrate our tests with three applications: happiness inequality in the US; self-assessed health in Europe; and sanitation ladders in Pakistan.

# Lay Summary

The first part of my thesis studies a problem in which a polluting firm can bribe an inspector to conceal evidence of illegal pollution. The government can deter bribes by rewarding the inspector to report the evidence, but doing so is costly. How can we reduce these costs?

Chapter 1, "Two Corruptible Monitors", asks whether more inspectors help to deter bribes. I compare different scenarios in which the government pays two inspectors to provide hard evidence about the firm's emissions. If the firm cannot bribe both monitors simultaneously, then the government might prefer to have each inspector receive a different piece of evidence, because then the firm can never conceal all of it. But if the firm can bribe both of the inspectors, then the government always prefers to have one inspector receive all the evidence. Why? Because when each inspector receives a different piece of evidence, the government must pay both inspectors in order to deter bribes to conceal one or other piece. But when one inspector receives both pieces of evidence, the same payment that deters bribes to conceal one piece, also deters bribes to conceal the other piece, so the government pays less overall. I conclude that (i) more monitors are not necessarily better, and (ii) it is important to consider the ease with which the firm can bribe multiple inspectors when designing incentives for more than one inspector.

Chapter 2, "Lemons by design: sowing secrets that curb corruption" with Andrew Clausen, asks whether the government can reduce costs by using random rewards and secrets to undermine trust between the firm and the inspector. We find that the regulator can do so by promising to pay a secret reward to one or the other whenever the inspector reports incriminating evidence; and by giving them each different clues about who stands to be rewarded. We find that the cheapest way to deter bribes is (i) to secretly select either the firm or the inspector to 'win' a reward whenever evidence is reported; and (ii) to give both the firm and the inspector a secret clue about who will win. If the inspector conceals the evidence, then the winner forgoes their reward, which means that they would have been better off not concealing the evidence. The clues are carefully constructed so that the firm only wants the inspector to conceal the evidence if his clue is informative enough to make him believe that the inspector is the winner. Similarly the inspector only wants to conceal evidence if her clue is informative enough to make her believe that the firm is the winner. If the firm anticipates that the inspector will only accept the bribe when she believes the firm will win, then it is rational for it to incorporate this into it's belief about the probably of winning, conditional on the inspector accepting the bribe. This causes the firm to require an even more informative clue in favour of not being rewarded, in order to induce them to accept a bribe. The inspector is undertaking the same reasoning about the firm's strategy, which causes her to require a

more informative clue in favour of not being rewarded, which in turn causes the firm to require a more informative clue, and so on. If the firm and the inspector keep reasoning in this way, then eventually neither of them will want to conceal the evidence, no matter what clue they receive. Therefore, they do not engage in bribery. As well as deterring bribes cheaply and robustly, this result proves that no matter how hard it is for two parties to agree to a deal in a situation where one of them knows something that there other does not, there always exists a situation where both of them know something that there other does not, for which it is strictly harder for them to agree.

The second part of my thesis studies the measurement of inequality in ordinal data such as self-reported wellbeing, sanitation, and perceived corruption. Chapter 3, "Inferring Inequality: Testing for Median-Preserving Spreads in Ordinal Data" with Ramses Abul Naga and Gaston Yalonetzky, devises a formal method for testing whether one data set is more unequal than another. We show that our tests are accurate in most applications. We illustrate our method with three applications: happiness inequality in the US; self-assessed health in Europe; and sanitation ladders in Pakistan.

# Acknowledgements

# Declaration

I declare that this thesis is an original report of my research, has been written by me, and has not been submitted for any previous degree. Chapter one is entirely my own work. Chapter two is a collaboration with Andrew Clausen. I am responsible for finding the optimal one-sided scheme, for writing a programme to help us find the optimal two-sided scheme and for writing the paper. Andrew directed our research towards random mechanisms, developed the coin toss example and encouraged me to develop my writing style. We both contributed to formalising the problem, finding the optimal two-sided scheme, and developing the proofs. Chapter three is a collaboration with Ramses Abul Naga and Gaston Yalonetzky. I am responsible for developing the theory, writing the programs, and conducting the Monte Carlo experiments. Ramses carried out the empirical investigations, and Gaston played a key role in developing the Z-statistic. We all contributed to writing and editing the paper. Due references have been provided on all supporting literatures and resources. All errors are my own responsibility.

Christopher Stapenhurst

# Contents

# Figures and Tables

## Figures

## Tables

# Chapter 1

# Two Corruptible Monitors

## Abstract

Do more monitors help to deter bribes? I compare optimal solutions for moral hazard problems in which a principal pays two corruptible monitors to provide hard evidence about an agent's action. If the agent can bribe one, but not both, of the monitors to conceal evidence, then the principal might prefer to have each monitor receive a different piece of evidence, because then the agent can never conceal all of it. But if the agent can bribe one or both of the monitors, then the principal always prefers to have one monitor receive all the evidence. Why? When each monitor receives a different piece of evidence, the principal must pay both monitors in order to deter bribes to conceal either piece. But when one monitor receives both pieces of evidence, the same payment that deters bribes to conceal one piece, also deters bribes to conceal the other piece, so the principal pays less overall. These results suggest that coalition formation frictions play an important role in incentive design problems with more than one monitor.

## 1.1   Introduction

The possibility of corruption is a serious problem in many incentive design problems. Whenever a principal relies on evidence provided by a third party monitor to reward or punish an agent, there is a risk that the agent might bribe the monitor to conceal the evidence. For example, factory owners bribe pollution inspectors to falsely report compliance (Duflo, Greenstone, Pande, & Ryan, 2013); construction firms bribe local politicians to grant planning permits[1]; firms bribe accountants to publish favourable audits[2]. Some institutions use multiple parties to monitor agents. For example, France now requires firms to be audited by two independent accountants (Vanstraelen, Richard, & R. Francis, 2009); the US False Claims Act rewards private citizens for providing evidence of Fraud to the Department of Justice (Beck,

---

1.   the Guardian. 2022. Spain's biggest corruption trial ends with 53 people convicted. [online] Available at: <https://www.theguardian.com/world/2013/oct/04/spain-corruption-trial-marbella-mayor-convicted> [Accessed 1 April 2022].
2.   Encyclopedia Britannica. 2022. Enron scandal | Summary, Explained, History, & Facts. [online] Available at: <https://www.britannica.com/event/Enron-scandal> [Accessed 1 April 2022].

2000); the US government allows private citizens to work independently of public agencies to provide evidence of environmental breeches (Langpap & Shimshack, 2010). One reason could be that more monitors can provide a large pool of information on which to condition incentive payments (e.g. Holmström, 1979). Another reason could be that competition between monitors may incentivize them to exert monitoring effort (e.g. Polinsky, 1980). A third possibility is that monitors may be able to deter each other from accepting bribes.

I focus on the third channel: do more independent monitors help to deter bribes? To answer this question, I study an incentive design problem in which a principal (she) designs incentives for an agent (he) to take the "right" action, by conditioning rewards for the agent on evidence provided by two monitors (they). Each monitor receives a binary signal. I interpret each signal as a piece of evidence that either materialises (outcome 1), or fails to materialise (outcome 0). For example, a piece of evidence could be a water sample taken to assess whether a nearby factory is emitting a legal amount of pollution. Evidence is informative and incriminating, because it is more likely to materialise when the agent takes the wrong action; but it is also noisy, because it might fail to materialise even though the agent takes the wrong action, or it might materialise when the agent takes the right action. It is without loss of generality to assume that one monitor's signal ("inferior" evidence) is noisier than the other's ("superior" evidence). Thus there are four possible outcomes no matter what action the agent chooses: no evidence, inferior evidence materialises, superior evidence materialises, both inferior and superior evidence materialise. In general, the principal incentivizes the agent to take the right action by rewarding exonerating outcomes, and/or punishing incriminating outcomes. A key result of this is paper will be that one of four agent-reward schemes is optimal. The "easy scheme" gives the agent a small reward whenever the monitors collectively report at most one piece of incriminating evidence; the "moderate scheme" gives the agent a moderate-sized reward so long as the monitors do not report superior evidence; and the "hard scheme" gives the agent a large reward only when neither monitor reports evidence. The fourth scheme, which I call the "cheeky" scheme, rewards the agent only when the monitors report the inferior evidence, but not the superior evidence. Most of the paper focuses on the easy, moderate and hard schemes both because they are more robust to bribes than the cheeky scheme, and because the cheeky scheme is only optimal for relatively extreme distributions of evidence.

Each monitor can conceal evidence by pretending that they didn't obtain it, when in fact they did. This possibility poses little threat to the principal in the absence of bribes, because the monitors are indifferent about reporting evidence (they do not face any costs of acquiring or reporting the evidence). But when the agent can bribe the monitors to conceal evidence, the principal must respond by designing rewards that induce the monitors to truthfully report their evidence. I assess the role of a second independent monitor by comparing three different scenarios that hold the total amount of available information fixed. In scenario 1, the benchmark case, one monitor accesses both pieces of evidence, and the agent can

bribe them to conceal one or both pieces. In scenario 2, each monitor accesses a different piece of evidence and the agent can bribe either one, but not both, monitors; In scenario 3, each monitor accesses only one type of evidence and the agent can bribe either one or both monitors.

In scenario 1, the agent will be willing to bribe the monitor in any outcome where the monitor can increase the agent's reward by concealing evidence. The principal can only deter bribes by paying the agent's reward to the monitor whenever the monitor reports evidence that forfeits the agent receiving it. Doing so ensures that the agent and the monitor receive the same joint surplus in all outcomes, so they can never increase their joint surplus by concealing evidence. Thus the principal can implement the easy scheme by rewarding the monitor whenever they report both pieces of evidence; she can implement the moderate scheme by rewarding the monitor whenever they report the superior evidence; and she can implement the hard scheme by rewarding the monitor whenever they report any type of evidence.

Which of the three schemes is the cheapest depends on the distribution of evidence. Two incriminating signals are incriminating, and two exonerating signals are exonerating; so all three schemes agree that the agent should be rewarded when both signals are exonerating, but not when both signals are incriminating. But what if one signal is exonerating and the other is incriminating? If the distribution of evidence is such that a single piece of exonerating evidence is enough to counteract a single piece of incriminating evidence then the principal should reward the agent whenever at least one signal is exonerating, so the easy scheme is the cheapest. This tends to be true if monitoring errors are more positively correlated conditional on the right action, than conditional on the wrong action. But if a single piece of incriminating evidence is enough to counteract a single piece of exonerating evidence then the principal should only reward the agent when both signals are exonerating, so the hard scheme is the cheapest. This tends to be true when monitoring errors are more positively correlated when the agent takes the wrong action. The remaining possibility is that the superior incriminating evidence is strong enough to counteract inferior exonerating evidence, but not vice versa. In this case, the moderate scheme is the cheapest. If monitoring errors are always perfectly correlated (i.e. if both monitors observe the same signal), then all three schemes cost the same.

In scenario 2, each monitor accesses a different type of evidence, and the agent can bribe one or other, but not both, of them to conceal it. The principal can implement the moderate reward scheme in the same way as for one monitor, because it only relies on the monitor with access to the superior type of evidence. The same is not true for the easy scheme. When both monitors have evidence, the agent can secure the reward by bribing either one of them. The principal can only deter such bribes by paying a reward to both of the monitors if and only if they both report evidence. This is more costly than the one monitor case, where the principal only had to pay one reward. On the other hand, the hard scheme is cheaper to implement

in this scenario than in scenario 1. In scenario 1, the principal needs to reward the monitor when they receive both pieces of evidence, because the monitor has the ability to conceal them both, thereby allowing the agent to claim the reward. But if the agent cannot coordinate with both monitors, then there is no need to pay either of them a reward when they both report evidence, because the agent cannot gain by bribing only one of them.

Whether the principal is better off in scenario 1 or scenario 2 depends on the distribution of evidence. If the distribution is one favouring the easy scheme (neither type of evidence is incriminating by itself), then the principal is better off with one monitor, because she can implement the easy scheme more cheaply than with two monitors. However, if the distribution is one favouring the hard scheme (both types of evidence are incriminating by themselves), then the principal is better off with two monitors, so long as the agent cannot bribe them both. In other cases, including all the cases where both monitors receive the same evidence, the principal is indifferent.

In scenario 3, each monitor accesses a different type of evidence, and the agent can bribe either one or both of them. The same easy and moderate schemes that deter bribes in scenario 2 continue to deter bribes here, because the agent only needs to conceal one piece of evidence to obtain a reward in those schemes. But hard reward scheme that deters bribes in scenario 2, no longer deters bribes when both monitors obtain evidence — the agent would like to conceal both pieces of evidence, and can easily bribe the monitors to do so since neither of them is being rewarded to report it. In order to make the hard scheme robust to coalitions containing both monitors, the principal must reward one or other (or a mix of both) monitors when they both report evidence. But then the total rewards paid to both monitors are exactly equal to the size of the reward paid in scenario 1. This implies that the principal is always weakly better off with a single monitor who accesses all the evidence, and strictly so when the distribution of evidence favours the easy reward scheme (i.e. when neither type of evidence is incriminating by itself). When both monitors access the same evidence, the easy and moderate schemes are weakly optimal, so once again, the principal is indifferent between the one and two monitors, even though the agent can bribe both of them.

I conclude that whether or not more monitors help to deter bribes depends on the ease with which the agent can bribe larger coalitions of monitors. If the agent cannot bribe both monitors simultaneously, and both types of evidence are sufficiently incriminating, then the principal is indeed better off with two, strategically independent monitors. If the agent can bribe both monitors then they are strategically equivalent to a single monitor. In this case, the principal is weakly better off with a single, omniscient monitor. In practice, the environment is likely to be somewhere between these two extremes because of frictions created by exogenous sources of private information, such as the precise characteristics of the evidence, the legal environment, or the players' 'moral' cost of bribery. It seems likely that larger coalitions suffer

more from these frictions, in which case it will be harder for the agent to bribe both monitors than one. Moreover, the principal may be able to endogenously create private information to undermine collusive coalitions (see chapter 2). In these cases, it will generally be preferable to have different monitors receive different pieces of evidence.

The problem of relying on a corruptible monitor to provide incentives for an agent taking a hidden action dates back to the seminal contribution of Tirole (1986), and has been followed by an extensive literature (e.g. Felli & Hortala-Vallve, 2016; Strausz, 1997; Vafaï, 2005). An extensive literature studies whether competition between monitors can incentivize costly monitoring effort (e.g. Liu, Wang, & Yin, 2022; McAfee, Mialon, & Mialon, 2008; Polinsky, 1980). Rahman (2012) tackles the same problem, but by having monitors monitor one another's effort, rather than by relying directly on market mechanisms. However, there are few results about the role that a multiplicity of monitors can play in deterring corruption. Kofman and Lawarrée (1993) introduce a second costly-but-incorruptible 'external' monitor. The optimal contract in their setting entails employing the internal monitor to monitor the agent and employing the external monitor to monitor the internal monitor. In all of these models, the principal has some means of obtaining her own unbiased information about the agent, either by observing output or by recourse to costly monitoring or by employing an external monitor. A key difference of my model is that the principal must rely exclusively on corruptible monitors to obtain information.

My problem is also closely related to the broader literature on mechanism design. Ben-Porath and Lipman (2012); J. R. Green and Laffont (1986); Koessler and Perez-Richet (2019), and others, study the role of hard evidence. Che and Kim (2006); Crémer (1996); J. Green and Laffont (1979); Laffont and Martimort (1997, 2000), and others, study the role of collusion on the set of social choice functions can be implemented. My principal also faces an implementation problem with hard evidence and collusion: she needs to implement a particular reward function for the agent (to induce him to take the right action), and she is able to choose transfers for the monitors in order to induce them to report their information. Hard evidence is strictly necessary for the principal because the monitors have no intrinsic preference over the agent's rewards, but they can always be bribed to report in the agent's favour.

Another related literature studies mutual monitoring in repeated games. Ben-Porath and Kahneman (1996) find that all individually rational payoffs can be attained with equilibrium strategies if and only if each player's action is observed by at least two other players. Aoyagi (2005) allow players to collude by means of a communication device, and find that individually rational payoffs can be attained with equilibrium strategies if and only if each player's action is almost perfectly observed by all other players collectively. I get a similar result in a static setting: when players can collude, the *number* of monitors observing the agent's action does not matter, only the *quality* of the information they collectively provide.

The outline of the chapter is as follows: Section 1.2 defines the players, their strategies and payoffs, and each of the three scenarios. Section 1.3 shows that the principal can restrict attention to collusion proof schemes. Section 1.4 solves the principal's problem in two stages. The first stage takes the agent's reward scheme as given and solves for the cheapest monitor wage scheme that deters bribes. The second stage then uses the costs implied by the monitor's wages to solve for the cheapest agent-reward scheme. Section 1.5 obtains the main result by comparing the optimal schemes across the three scenarios. Section 1.6 discusses the robustness of the baseline model by considering different assumptions about limited liability, collusion, and the distribution of evidence. Section 1.7 discusses the results in the context of applications to journalism, whistle-blowing and financial auditing. Section 1.8 concludes.

## 1.2 Model

A risk neutral agent (he) chooses between a "right" action and a "wrong" action. Other things equal, the right action yields a payoff of 0, whereas the wrong action yields a payoff of 1. There are two binary signals $s_1$ and $s_2$ which are informative about the agent's action. Each signal takes the values 1 ("incriminating evidence") or 0 ("no evidence"), so the joint signal $s = (s_1, s_2)$ is an element of the set $S := \{0,1\}^2$ (where ":=" denotes the definition of a symbol). The distribution of $s$ depends on the agent's choice of action. If the agent takes the wrong action, then $s$ occurs with probability $\tau_s$; if the agent takes the right action, then $s$ occurs with probability $\pi_s$. If $\Delta_s := \pi_s - \tau_s > 0$, then $s$ is more likely to occur when the agent takes the right action than when he takes the wrong action, so we can think of $s$ as *exonerating evidence*. Similarly, if $\Delta_s < 0$, then $s$ constitutes *incriminating evidence*, and if $\Delta_s = 0$, then it is neutral. Note that $\sum_{s \in S} \pi_s = \sum_{s \in S} \tau_s = 1$ implies $\sum_{s \in S} \Delta_s = 0$, so a realisation $s$ can only be incriminating if some other realisation, $s'$, is exonerating. The main results assume that $\Delta_{11} < 0$ and $\Delta_{00} > 0$, so two pieces of incriminating evidence are indeed incriminating, and two pieces of exonerating evidence are indeed exonerating. Moreover, I assume that $\Delta_{11} \leq \Delta_{10} \leq \Delta_{00}$ and $\Delta_{11} \leq \Delta_{01} \leq \Delta_{00}$, so a single piece of evidence is more incriminating than no evidence, but less incriminating than both pieces of evidence. These assumptions imply that $\Delta_{00} + \Delta_{01} > \Delta_{10} + \Delta_{11}$ and $\Delta_{00} + \Delta_{10} > \Delta_{01} + \Delta_{11}$ which means that $s_i = 1$ is incriminating when $s_{-i}$ is not observed for $i = 1, 2$. It is without loss of generality to further assume that $\Delta_{01} \geq \Delta_{10}$. Consequently, $\Delta_{01} + \Delta_{11} \geq \Delta_{10} + \Delta_{11}$ and $\Delta_{00} + \Delta_{01} \geq \Delta_{00} + \Delta_{10}$, which means that $s_1 = 1$ is more incriminating than $s_2 = 1$, and $s_1 = 0$ is more exonerating than $s_2 = 0$. Thus the (superior) signal $s_1$ is more accurate than the (inferior) signal $s_2$. If $\Delta_{00} = 1$ (resp. $\Delta_{11} = -1$) then we arrive in the case of perfect monitoring because it must be that $\pi_{00} = 1$ and $\tau_{00} = 0$ (resp. $\pi_{00} = 0$ and $\tau_{00} = 1$). In this special case, the principal can perfectly infer that the agent must have taken the right action if and only if $s = (0,0)$ (resp. $s \neq (1,1)$). Other cases are considered in section 1.6.

The principal wants to incentivize the agent to take the right action. She can only do so by committing to reward or punish him on the basis of the signal $s$. Specifically, the principal promises to pay the agent a transfer $T(s)$ when signal $s$ is reported. This transfer function $T : S \to \mathbb{R}$ needs to satisfy

$$\sum_{s \in S} \pi_s T(s) \geq 1 + \sum_{s \in S} \tau_s T(s)$$

or equivalently,

$$\sum_{s \in S} \Delta_s T(s) \geq 1, \tag{IC}$$

to ensure that the agent's payoff from taking the right action is greater than his payoff from taking the wrong action. This is the agent's *incentive compatibility* or (IC) constraint. I assume that indifferences are broken in favour of the principal, so the agent complies if and only if (IC) is satisfied. It is clear that this constraint must be satisfied by either rewarding the agent when exonerating evidence is realised, or by punishing the agent when incriminating evidence is realised, because these configurations will ensure that the product $\Delta_s T(s)$ is positive. For most of the paper, I rule out punishments by assuming that the agent has limited liability, so the principal must rely on rewarding exonerating evidence. If (IC) is satisfied, then the agent chooses the right action in equilibrium so the cost of transfer function $T$ is $\sum_{s \in S} \pi_s T(s)$.

If the principal directly observed both signals $s_1$ and $s_2$, then her objective would simply be to choose a transfer function to minimise the cost of satisfying the agent's (IC) constraint. However, the key feature of the model is that the signal is not observed by the principal, but by either one or two monitors. Moreover, the monitors are corruptible: they can be bribed to conceal evidence. If monitor $i \in \{1, 2\}$ receives a 1 signal, then they may send either a 1 or a 0 report to the principal. But if they receive a 0 signal, then they have no choice but to send a 0 report to the principal. The interpretation is that 1 is hard incriminating evidence, and 0 is a dearth of hard incriminating evidence. Hard evidence can be concealed, but not fabricated. Define a component-wise partial ordering, '$\leq$', by $m \leq s$ if and only if $m_i \leq s_i$ for all $i = 1, 2$, which is true only if $m$ can be obtained from $s$ by fully or partially concealing evidence. The 'lower contour set' of an outcome $s \in S$ is the set of outcomes that can be reached from it by concealing evidence, $\{m \in S \mid m_i \leq s_i, i = 1, 2\}$.

The principal anticipates the possibility that the agent could bribe the monitor to conceal evidence. For instance, if the principal chooses $T(0, 0) > T(1, 0)$, then the agent would be willing to bribe the monitor anything up to $T(0, 0) - T(1, 0)$ to report $(0, 0)$ instead of $(1, 0)$. The principal can mitigate against bribery by rewarding the monitors for reporting hard evidence. It does so by committing to a wage function $w_i : S \to \mathbb{R}_+$ which specifies how much monitor $i$ gets paid when they collectively report evidence $s$. The main result assumes

that both monitors have limited liability, so wages must be weakly positive (this assumption is relaxed in section 1.6). A *scheme* $(T, w_1, w_2)$ specifies a transfer function for the agent together with a pair of wage functions for the monitors. The cost of a scheme $(T, w_1, w_2)$ is $c(T, w_1, w_2) := \sum_{s \in S} \pi(s)(T(s) + w_1(s) + w_2(s))$. The principal's problem is to choose a scheme of transfers and wages to minimise the expected cost, subject to the agent's incentive compatibility constraint, the agent's liability constraint, the monitors' liability constraints.

The agent has an incentive to bribe the monitors to conceal evidence whenever doing so would lead to him being rewarded. I assume that bribery negotiations take place after the monitor obtains evidence (if the monitor does not obtain evidence then there is no scope for bribes because the monitor cannot change the monitor's reward), and that the evidence realisation $s$ is common knowledge. The latter assumption ensures that any scheme that is robust to bribes in this environment, will continue to be robust to bribes in environments where the players negotiate bribes without perfect knowledge of the evidence realisation. It does not necessarily conflict with the principal's need to rely on monitors to report the evidence, because evidence needs to be verifiable in most practical contexts. Here, I restrict attention to deterministic bribes — all the results continue to hold when stochastic bribes are permitted.

The precise form of a bribe depends on which monitor(s) have evidence and which collusive coalitions can form. I compare three scenarios:

- **Scenario 1.** Monitor 1 receives both signals $s_1$ and $s_2$, and the agent can bribe the monitor to conceal one or both of them. In outcome $s$, a bribe $(b(s), m(s))$ specifies a report $m(s) = (m_1(s), m_2(s))$ and with $m \leq s$ for $i = 1, 2$, and a bribe $b$ to be paid to monitor 1. The agent is willing to pay the (non-trivial) bribe if

$$T(m(s)) - b(s) > T(s), \tag{1.1}$$

  and monitor 1 is willing to accept the (non-trivial) bribe if

$$w_1(m(s)) + b(s) > w_1(s). \tag{1.2}$$

- **Scenario 2.** Monitor $i$ receives signal $s_i$, and the agent can either bribe one, but not both of them. In outcome $s$, a bribe $(b_i(s), m_i(s))$ specifies a report $m_i(s) \leq s_i$, and a bribe $b_i$ to be paid to monitor $i = 1, 2$. The agent is willing to pay the bribe if

$$T(m_i(s), s_{-i}) - b_i(s) > T(s), \tag{1.3}$$

  and monitor $i$ is willing to accept the bribe if

$$w_i(m_i(s), s_{-i}) + b_i(s) > w_i(s). \tag{1.4}$$

- **Scenario 3.** Monitor $i$ receives signal $s_i$, and the agent can bribe either one or both of them. In outcome $s$, a bribe $(b_1(s), b_2(s), m_1(s), m_2(s))$ specifies reports $m_i(s) \leq s_i$, and bribes $b_i$, to be paid to one or both monitors $i = 1, 2$. The agent is willing to bribe one monitor if (1.3) holds, or both monitors if

$$T(m_1(s), m_2(s)) - \sum_{i=1,2} b_i(s) > T(s) \tag{1.5}$$

  holds; and monitor $i = 1, 2$ is willing to accept the bribe if

$$w_i(m_1(s), m_2(s)) + b_i(s) > w_i(s). \tag{1.6}$$

The inequalities are strict because of my assumption that indifferences are broken in favour of the principal. This assumption simplifies the problem because it will ensure that the principal's feasible set is compact; but has no substantial bearing on the solution because the principal can give the players a strict preference to reject bribes by adding some arbitrarily small amount to the monitors' wages. A bribe is *feasible* in outcome $s$ if the monitor is willing to accept it and the agent is willing to pay it. I discuss the implications of allowing other coalitions (excluding the agent and/or including the principal) in section 1.6. *How* the players choose between different feasible bribes does not matter, because lemma 2 will show that the principal can restrict attention to schemes for which there are no feasible, non-trivial bribes.

How do bribes affect the principal's problem? The principal must use a 'robustly incentive compatible' scheme in oder to incentivize the agent to take the right action. A scheme is *robustly incentive compatible* if it satisfies (IC), no matter what feasible bribes take place. For example, in scenario 3, a scheme $(T, w_1, w_2)$ must satisfy

$$\sum_{s \in S} \Delta_s \left( T(m_1(s), m_2(s)) - b_1(s) - b_2(s) \right) \geq 1, \tag{RIC}$$

for every feasible bribe $(b(s), m(s))$, and every outcome $s \in S$. Bribes also affect the cost of the scheme. Erring on the side of caution, I assume that the cost of a scheme is equal to the maximum amount that the principal might have to pay out in rewards following any feasible bribes. Indeed, this is how much the principal pays if the agent and the monitor conceal evidence to maximise their joint surplus. For example, in scenario 1, the cost of a scheme $(T, w_1, w_2)$ is

$$\sup_{\substack{(b(s), m(s)) \\ \text{is feasible}}} \sum_{s \in S} \pi_s (T(m(s)) + w_1(m(s)) + w_2(m(s))).$$

The cost is defined analogously in other scenarios, with the appropriate bribe as a choice variable.

I conclude this section by summarising the timing of the model:

    i) the principal proposes a scheme;

   ii) the agent chooses either the right or the wrong action;

  iii) the monitor(s) obtain evidence $s$;

  iv) the players negotiate a bribe;

   v) the monitor(s) report evidence;

  vi) all players' payoffs are realised.

## 1.3 Bribe Proof Schemes

Anticipating all the possible bribes that may take place is unnecessary, because the principal can simplify her problem by restricting attention to "bribe proof" schemes in which all bribes are deterred. A scheme $(T, w_1, w_2)$ is bribe proof if and only if there are no feasible bribes in any outcome $s$. Precisely what this means will depend on which scenario is under consideration, because the scenario determines which bribes are feasible. Specifically, a scheme is *j-bribe proof* if it is bribe proof in scenario $j = 1, 2, 3$. Lemma 1 shows that a scheme is $j$-bribe proof if and only if no coalition can generate surplus by suppressing evidence.

**Lemma 1.** *A scheme $(T, w_1, w_2)$ is*

- *1-bribe proof if and only if*

$$T(s) + w_1(s) \geq T(m) + w_1(m) \tag{1-BP}$$

    *for all $m \leq s$;*

- *2-bribe proof if and only if*

$$T(s) + w_i(s) \geq T(m_i, s_{-i}) + w_i(m_i, s_{-i}) \tag{2-BP}$$

    *for all $m_i \leq s_i$, and for all $i = 1, 2$;*

- *3-bribe proof if and only if it is 2-bribe and*

$$T(s) + w_1(s) + w_2(s) \geq T(m) + w_1(m) + w_2(m) \tag{3-BP}$$

    *for all $m \leq s$.*

*Proof.* I first prove the 'if' statement for each scenario, and then the 'only if' statement.

Let $(T, w_1, w_2)$ be a scheme, and let $(b(s), m(s))$ be a feasible bribe in some outcome $s$ in scenario 1, with $m \neq s$. Then $(b(s), m(s))$ must strictly satisfy inequalities (1.1) and (1.2). Summing these gives $w_1(m(s)) + T(m(s)) > w_1(s) + T(s)$, which is the opposite of (1-BP). Therefore if (1-BP) holds then the scheme is 1-bribe proof, proving the "if" statement for scenario 1. For scenario 2, let $(b_i(s), m_i(s))$ be a feasible bribe with $m_i \neq s_i$. Then $(b_i(s), m_i(s))$

must satisfy inequalities (1.3) and (1.4), which sum to give $T(m_i(s), s_i) + w_i(m_i(s), s_{-i}) > T(s) + w_i(s)$. This is ruled out by (2-BP), proving the "if" statement for scenario 2. Finally, let $(b_1(s), b_2(s), m_1(s), m_2(s))$ be a feasible bribe in scenario 3, with either $m_1 \neq s_1$ or $m_2 \neq s_2$. Then $(b_1(s), b_2(s), m_1(s), m_2(s))$ either satisfies (1.3) and (1.4) for some monitor $i$, in which case it is ruled out by (2-BP); or else it satisfies inequalities (1.5) and (1.4) for monitors $i = 1, 2$. The latter inequalities sum to give $T(m_1(s), m_2(s)) + \sum_{i=1,2} w_i(m_1(s), m_2(s)) \geq T(s) + \sum_{i=1,2} w_i(s)$, which is ruled out by (3-BP). This proves the "if" statement for scenario 3.

Now suppose that one of the constraints is violated. Then some coalition of players can generate surplus by concealing evidence. Any bribe that splits this surplus between the members of the coalition will be feasible. For example, if (3-BP) is violated in a state $s$ with report $m$, then there exists some $\epsilon > 0$ such that

$$T(m_1(s), m_2(s)) - \sum_{i=1,2} (w_i(m_1(s), m_2(s)) - w_i(s)) > T(s) + 2\epsilon$$

and $w_i(m_1(s), m_2(s)) + w_i(s) - w_i(m_1(s), m_2(s)) + \epsilon > w_i(s)$, so the bribes $b_i(s) = w_i(s) - w_i(m_1(s), m_2(s)) + \epsilon$ for monitors $i = 1, 2$ are feasible. Therefore $(T, w_1, w_2)$ is not 3-bribe proof, proving the "only if" part of the statement. The proofs for scenarios 1 and 2 are analogous. $\square$

The scenario 2 constraints closely resemble the scenario 1 constraints, except with the addition of the $i$ subscripts, and the lack of any constraint on transfers and wages in outcome (1,1) relative to outcome (0,0). This is because the agent cannot bribe both monitors, and so the principal does not have to worry about both pieces of evidence being concealed at the same time. In scenario 3, all but one of these constraints involves the agent bribing both monitors, even though only one monitor suppresses their evidence. It is necessary to account for such coalitions to rule out schemes that increase one monitor's wage when the other monitor reports a 0. In such schemes, the agent benefits from including the second monitor in the coalition because doing so increases the overall surplus that the coalition generates by concealing just one signal.

Lemma 2 shows that, in scenarios 1 and 3, the principal can restrict attention to schemes for which there are no (non-trivial) feasible bribes.

**Lemma 2.** *If a scheme $(T, w_1, w_2)$ respects robust incentive compatibility and limited liability, then there exist 1- and 3-bribe proof schemes $(T^1, w^1)$ and $(T^3, w^3)$ that respect incentive compatibility and limited liability, and have weakly lower cost, $c(T^j, w^j) \leq c(T, w)$ for $j = 1, 3$.*

*Proof.* The proof of lemma 2 follows a revelation principle type argument (Laffont & Martimort, 1997). I give the proof for scenario 3; the proof for scenario 1 is analogous. Suppose a scheme $(T, w_1, w_2)$ satisfies (RIC), limited liability, and costs $C$, but violates 3-bribe proofness. Lemma 1 shows that there exists a feasible bribe for each state $s$ where (3-BP) is violated.

The set of feasible bribes is not compact (because of the strict inequalities in equations (1.1)–(1.6)), and therefore there may not exist a feasible bribe that maximises the agent's payoff. But the agent's payoff is bound above by $\max_s T(s) + w_1(s) + w_2(s)$, so it has a supremum, and this supremum is attained by a bribe $(b_1(s), b_2(s), m_1(s), m_2(s))$ in the closure of the feasible set.

Define a new scheme by $T^3(s) = T(m_1(s), m_2(s)) - b_1(s) - b_2(s)$ and $w_i^3(s) = w_i(m_1(s), m_2(s)) + b_i(s)$. Then $(T^3, w_1^3, w_2^3)$ has the same ex post payoffs as the old scheme $(T, w_1, w_2)$ together with the bribe. It satisfies limited liability because $(T, w_1, w_2)$ satisfies limited liability and $(b_1(s), b_2(s), m_1(s), m_2(s))$ is in the closure of the set of feasible bribes, so $T(m_1(s), m_2(s)) - b_1(s) - b_2(s) \geq T(s) \geq 0$ and $w_i(m_1(s), m_2(s)) + b_i(s) \geq w_i \geq 0$. Similarly, $(T^3, w_1^3, w_2^3)$ satisfies incentive compatibility because $(T, w_1, w_2)$ satisfies (RIC) so $\sum_{s \in S} \pi_s T^3(s) = \sum_{s \in S} \pi_s T(s)(T(m_1(s), m_2(s)) - b_1(s) - b_2(s)) \geq 1$. Moreover, $(T^3, w_1^3, w_2^3)$ must be bribe proof. If it were not, then there would be a bribe that strictly increases the agent's payoff $(b_1'(s), b_2'(s), m_1'(s), m_2'(s))$. But then the composition of the two bribes, $(b_1(s) + b_1'(s), b_2(s) + b_2'(s), m_1'(m_1(s)), m_2'(m_2(s)))$, is feasible under $(T, w_1, w_2)$ and yields the agent a higher payoff than $(b_1(s), b_2(s), m_1(s), m_2(s))$, contradicting the fact that $(b_1(s), b_2(s), m_1(s), m_2(s))$ was assumed to attain the supremum of the agent's payoffs across all feasible bribes. Finally, the scheme $(T^3, w_1^3, w_2^3)$ has a weakly lower cost because

$$
\begin{aligned}
c(T, w_1, w_2) &= \max_{\substack{(b_1(s), b_2(s), m_1(s), m_2(s)) \\ \text{is feasible}}} \sum_{s \in S} \pi_s(T(m'(s)) + w_1(m'(s)) + w_2(m'(s))) \\
&\geq \sum_{s \in S} \pi_s(T(m(s)) + w_1(m(s)) + w_2(m(s))) \\
&= \sum_{s \in S} \pi_s(T^3(s) + w_1^3(s) + w_2^3(s)) \\
&= c(T^3, w_1^3, w_2^3),
\end{aligned}
$$

where the second equality holds because the bribes cancel out. □

The proof of lemma 2 does not work for scenario 2, because the composition of two bribes is not feasible if they each involve different monitors. But the purpose of studying scenario 2 is to show that the principal can potentially save money by hiring two monitors if they cannot be bribed simultaneously. If lemma 2 does not hold for scenario 2, then it only means that the cheapest feasible, bribe proof solution gives an upper bound on the cheapest feasible solution.

Lemmas 1 and 2 imply that the principal's problem in each scenario can be formulated as follows:

$$\max_{T, w_1, w2} -\sum_{s \in S} \pi_s (T(s) + w_1(s) + w_2(s))$$

$$\text{st.} \sum_{s \in S} \Delta_s T(s) \geq 1 \tag{IC}$$

$$T(s) \geq 0 \qquad \qquad \forall s \in S \tag{LLA}$$

$$w_i(s) \geq 0 \qquad \qquad \forall s \in S, i = 1, 2 \tag{LLM}$$

$$(j - \text{BP})$$

where $j = 1$ (in scenario 1), $j = 2$ (in scenario 2) or $j = 2, 3$ (in scenario 3).

## 1.4 Optimal schemes

I solve for the optimal schemes in two steps. First I take the transfer function $T$ as given, and solve for the cheapest wages $w_1$ and $w_2$ that deter bribes in each scenario. Then, given the cost of the transfer $T$ and the corresponding wages $w_1$ and $w_2$, I find the cheapest scheme for given parameters $\pi$ and $\tau$.

### 1.4.1 incentivizing the monitors

Proposition 1 says that the monitors are collectively paid the marginal value of their evidence in any optimal scheme. The only differences between the three scenarios occur in the outcome $(1, 1)$, because it is the only outcome where the marginal value of a piece of evidence depends on whether the other evidence is reported.

**Proposition 1.** *If $(T^j, w_1^j, w_2^j)$ is a solution to the principal's problem in scenario $j$, then*

$$w_1^j(0, 0) + w_2^j(0, 0) = W^j(0, 0; T^j) := 0 \tag{1.7}$$

$$w_1^j(1, 0) + w_2^j(1, 0) = W^j(1, 0; T^j) := \max\{0, T^j(0, 0) - T^j(1, 0)\} \tag{1.8}$$

$$w_1^j(0, 1) + w_2^j(0, 1) = W^j(0, 1; T^j) := \max\{0, T^j(0, 0) - T^j(0, 1)\} \tag{1.9}$$

*for $j = 1, 2, 3$, and*

$$\begin{aligned} w_1^1(1, 1) + w_2^1(1, 1) = W^1(1, 1; T^j) := \max\{0, T^j(0, 1) - T^j(1, 1), T^j(1, 0) - T^j(1, 1), \\ T^j(0, 0) - T^j(1, 1)\} \end{aligned} \tag{1.10}$$

$$\begin{aligned} w_1^2(1, 1) + w_2^2(1, 1) = W^2(1, 1; T^j) := \max\{0, T^j(0, 1) - T^j(1, 1), T^j(1, 0) - T^j(1, 1), \\ T^j(0, 1) + T^j(1, 0) - 2T^j(1, 1)\} \end{aligned} \tag{1.11}$$

$$w_1^3(1, 1) + w_2^3(1, 1) = W^3(1, 1; T^j) := \max\{w_1^1(1, 1) + w_2^1(1, 1), w_1^2(1, 1) + w_2^2(1, 1)\}. \tag{1.12}$$

*Proof.* Equality (1.7) says that neither monitor gets paid if neither has evidence. The wages $w_1(0,0)$ and $w_2(0,0)$ feature negatively in the principal's objective and no coalition can ever benefit from monitor $i$ concealing evidence in these outcomes. So the corresponding wages only feature positively in monitor $i$'s limited liability constraints, which means that these constraints must hold with equality in any optimal solution, giving (1.7).

Similarly, $w_1(1,0)$ and $w_2(1,0)$ have negative coefficients in the principal's objective and in some bribe constraints. The only places where they have positive coefficients are in the liability constraints, and in the bribe constraints for the outcome $(1,0)$. Substituting (1.7) into (1-BP) for outcome $(1,0)$ gives $w_1(1,0) \geq T(0,0) - T(1,0)$. The limited liability constraints give $w_1(1,0) \geq 0$ and $w_2(1,0) \geq 0$, so any optimal solution must have $w_2(1,0) = 0$ and $w_1(1,0) = \max\{0, T(0,0) - T(1,0)\}$, hence equation (1.8) holds in scenario 1. The same argument applies for scenario 2. Scenario 3, has an additional bribe constraint, (3-BP), which says that $w_1(1,0) + w_2(1,0) \geq T(0,0) - T(1,0)$, but this is already satisfied by (1.8). The same argument applies for outcome $(0,1)$.

The optimal wages differ between scenarios when both monitors have evidence. In scenario 1, $w_1(1,1)$ enters three bribe constraints, one for each of the possible deviations to $(1,0), (0,1)$ and $(0,0)$. Substituting $w_1(1,0) = \max\{0, T(0,0) - T(1,0)\}$, $w_1(0,1) = \max\{0, T(0,0) - T(0,1)\}$ and $w_1(0,0) = 0$ into these constraints gives

$$w_1(1,1) \geq T(1,0) - T(1,1) + \max\{0, T(0,0) - T(1,0)\}$$
$$w_1(1,1) \geq T(0,1) - T(1,1) + \max\{0, T(0,0) - T(0,1)\}$$
$$w_1(1,1) \geq T(0,0) - T(1,1).$$

The upper envelope of these constraints gives $w_1(1,1) \geq \max\{0, T(0,0) - T(1,1), T(1,0) - T(1,1), T(0,1) - T(1,1)\}$. The wage $w_2(1,1)$ does not enter any bribe constraints but limited liability requires $w_2(1,1) \geq 0$. Both feature negatively in the principal's objective, so both must hold with equality. Summing gives (1.10).

In scenario 2, $w_1(1,1)$ only enters one bribe constraint, which says that $w_1(1,1) \geq T(0,1) - T(1,1)$, as well as monitor 1's liability constraint. Wage $w_1(1,1)$ has a negative coefficient in the objective, so one of these two constraints must hold with equality in any optimal solution, so $w_1^2(1,1) = \max\{0, T(0,1) - T(1,1)\}$. The same argument applies for $w_2^2(1,1)$. Summing gives (1.11).

In scenario 3, $w_1(1,1)$ and $w_2(1,1)$ enter five bribe constraints. Substituting (1.8) and (1.9) into them yields

$$w_1(1,1) \geq T(0,1) - T(1,1)$$
$$w_2(1,1) \geq T(0,1) - T(1,1)$$
$$w_1(1,1) + w_2(1,1) \geq T(0,1) - T(1,1) + \max\{0, T(0,0) - T(0,1)\}$$
$$w_1(1,1) + w_2(1,1) \geq T(1,0) - T(1,1) + \max\{0, T(0,0) - T(1,0)\}$$
$$w_1(1,1) + w_2(1,1) \geq T(0,0) - T(1,1),$$

as well as their liability constraints, which together imply that $w_1(1,1) + w_2(1,1) \geq 0$. The first two bribe constraints sum to give $w_1(1,1) + w_2(1,1) \geq T(0,1) + T(0,1) - 2T(1,1)$. The objective is decreasing in $w_1(1,1) + w_2(1,1)$, so at least one of these five constraints must hold with equality. If it is not the case that $w_1^3(1,1) + w_2^3(1,1) = T(0,0) - T(1,1)$, then $w_1^3(1,1) + w_2^3(1,1) = w_1^2(1,1) + w_2^2(1,1) > w_1^1(1,1) + w_2^1(1,1)$, so (1.12) holds. Otherwise, if $w_1^3(1,1) + w_2^3(1,1) = T(0,0) - T(1,1)$, then $w_1^3(1,1) + w_2^3(1,1) = w_1^1(1,1) + w_2^1(1,1) > w_1^2(1,1) + w_2^2(1,1)$, and (1.12) still holds. □

We can already see from equation (1.12) that the value of a second monitor is limited, because the principal is more constrained in scenario 3 than in scenario 1. Proposition 2 will show how this difference in the wages will affect the optimal choice of transfer.

The function $W^j(s; T)$ defined by proposition 1 allows us to remove the bribe constraints and the monitors' liability constraints from the principal's problem. Substituting it gives

$$\max_T - \sum_{s \in S} \pi_s (T(s) + W^j(s; T))$$
$$\text{st. } \sum_{s \in S} \Delta_s T(s) \geq 1 \qquad\qquad\qquad \text{(IC)}$$
$$T(s) \geq 0 \qquad\qquad\qquad\qquad \forall s \in S. \qquad \text{(LLA)}$$

### 1.4.2   incentivizing the agent

The principal incentivizes the agent by rewarding him for exonerating outcomes, i.e. outcomes $s$ with $\Delta_s > 0$. She trades this off against the cost of rewarding to the agent, which is determined by the probability of realising the outcome $s$ when the agent takes the right action; and the cost of deterring the monitors from falsely reporting the outcome $s$, which is determined by the probability of realising an outcome $s'$ from which the agent can bribe the monitors to report $s$.

**Lemma 3.** *The optimal transfers for the agent are bang-bang. If $(T^*, w_1^*, w_2^*)$ is a solution to the principal's problem, then there exists a set of 'reward' outcomes $R \subset S$ such that the agent is given a reward of size $\frac{1}{\sum_{s \in R} \Delta_s}$ if the monitors report an outcome in $R$. In all other outcomes, he receives no reward. I.e. $T^*(s) = \mathcal{I}(s \in R) / \sum_{s \in R} \Delta_s$.*

*Proof.* I first show that the solution to the principal's problem coincides with a solution to a linear programme, and then show that the solution to this linear programme must be bang-bang.

Any solution $(T^*, w_1^*, w_2^*)$ to any of the three problems implies an ordering over rewards, e.g. $T^*(0,1) \geq T^*(1,0) \geq T^*(0,0) \geq T^*(1,1)$. This means that $(T^*, w_1^*, w_2^*)$ also solves the problem where these three inequalities are imposed as constraints. This *order-constrained* problem has four choice variables, $T(0,0), T(1,0), T(0,1)$ and $T(1,1)$, and five linear constraints (the (IC) constraint, the three order constraints, and one active liability constraint). The role of the active liability constraint ($T(1,1) \geq 0$ in the example) is the same as the role of the order constraints, so I will refer to the order constraints and the liability constraint collectively as 'order constraints' for the remainder of the proof. The ordering determines which arms of the maximands in (1.8)–(1.12) are active, so the objective of the order constrained problem is linear. For example, if $T(0,1) \geq T(1,0) \geq T(0,0) \geq T(1,1)$, then proposition 1 implies that scenario 1 has

$$w_1^1(0,0) + w_2^1(0,0) = 0$$
$$w_1^1(1,0) + w_2^1(1,0) = 0$$
$$w_1^1(0,1) + w_2^1(0,1) = 0$$
$$w_1^1(1,1) + w_2^1(1,1) = T(1,0) - T(1,1),$$

so the scenario 1 problem reduces to the linear programme

$$\min_T \pi_{00} T(0,0) + (\pi_{10} + \pi_{11}) T(1,0) + \pi_{01} T(0,1)$$

$$\text{s.t. } \sum_{s \in S} \Delta_s T(s) \geq 1, \tag{IC}$$

$$T(0,1) \geq T(1,0)$$

$$T(1,0) \geq T(0,0)$$

$$T(0,0) \geq T(1,1)$$

$$T(1,1) \geq 0. \tag{LLA}$$

We can solve this problem by studying the Khun-Tucker conditions. The first order conditions yield a system of four equations (one for each choice variable) in five Lagrange multipliers (one for each constraint). For example, if $\lambda_{IC}$ denotes the multiplier on (IC) and $\lambda_s$ denotes the multiplier on the order constraint for $T(s)$, then we get

$$\lambda_{IC}\Delta_{11} + \lambda_{11} - \lambda_{00} = 0$$
$$-\pi_{00} + \lambda_{IC}\Delta_{00} + \lambda_{00} - \lambda_{10} = 0$$
$$-\pi_{10} - \pi_{11} + \lambda_{IC}\Delta_{10} + \lambda_{10} - \lambda_{01} = 0$$
$$-\pi_{01} + \lambda_{IC}\Delta_{01} + \lambda_{01} = 0$$

If two multipliers equal zero then we have a system of three equations in two unknowns, which only has a solution in certain, knife-edge cases.[3] Therefore, at most one of the multipliers can be generically equal to 0. On the other hand, if the four order constraint multipliers are all positive, then complementary slackness implies that these constraints all hold with equality, so the transfers would all equal 0, violating (IC). Therefore, at least one of the order multipliers $\lambda_s$ must equal 0.

Together, these conclusions imply that exactly one of the order multipliers equals zero, so complementary slackness implies that exactly one of the corresponding inequalities is strict. This implies that each of the transfers is either (i) equal to zero, if it lies below the strict inequality; or it is (ii) equal to the largest transfer, if it lies above the strict inequality. The set $R$ is equal to the set of outcomes with transfers above the strict inequality. In the example above, if $\lambda_{10} = 0$, then we get $T(0,1) = T(1,0) > T(0,0) = 0$, so $R = \{(0,1),(1,0)\}$.

Let $T^*$ denote the size of the strictly positive transfer. The principal wants to minimise $T^*$ subject to an IC constraint which is increasing in $T^*$. Therefore the IC constraint must hold with equality in any optimal solution, giving $\sum_{s \in R} \Delta_s T^* = 1$, or $T^* = 1/\sum_{s \in R} \Delta_s$. □

---

3. These cases admit a continuum of solutions, at least two of them will have the binary structure described in the statement of lemma 3.

Proposition 2 formally describes the easy, moderate and hard schemes described in the introduction. To aid intuition, I present these an matrices where entry $(i, j)$ in the wage matrix denotes the value of the corresponding function in outcome $(i - 1, j - 1)$. For example, the easy, moderate and hard reward functions are given by

$$T_E := \begin{bmatrix} \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} & \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} \\ \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} & 0 \end{bmatrix}$$

$$T_M := \begin{bmatrix} \frac{1}{\Delta_{00}+\Delta_{01}} & \frac{1}{\Delta_{00}+\Delta_{01}} \\ 0 & 0 \end{bmatrix}$$

$$T_H := \begin{bmatrix} \frac{1}{\Delta_{00}} & 0 \\ 0 & 0 \end{bmatrix}.$$

There is also a fourth type of scheme that may be optimal, which I refer to as the *cheeky* scheme with transfer function

$$T_C := \begin{bmatrix} 0 & \frac{1}{\Delta_{01}} \\ 0 & 0 \end{bmatrix}.$$

In this scheme, the principal gains by reducing the probability of paying any rewards by only paying rewards when monitor 2 reports incriminating evidence, which could be very rarely, given that the agent takes the right action in equilibrium. In section 1.6 I argue that this scheme is less robust than the others because it is susceptible to corruption by the principal.

**Proposition 2.** *One of the following four schemes is optimal in scenario 1 (monitor 1 receives both signals $s_1$ and $s_2$):*

$$S_E^1 := \left( T_E, \begin{bmatrix} 0 & 0 \\ 0 & \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \right)$$

$$S_M := \left( T_M, \begin{bmatrix} 0 & 0 \\ \frac{1}{\Delta_{00}+\Delta_{01}} & \frac{1}{\Delta_{00}+\Delta_{01}} \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \right)$$

$$S_H^1 := \left( T_H, \begin{bmatrix} 0 & \frac{1}{\Delta_{00}} \\ \frac{1}{\Delta_{00}} & \frac{1}{\Delta_{00}} \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \right)$$

$$S_C := \left( T_C, \begin{bmatrix} 0 & 0 \\ 0 & \frac{1}{\Delta_{01}} \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \right).$$

*The costs of these schemes are $c_E^1(\pi, \Delta) := \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}}$, $c_M(\pi, \Delta) := \frac{1}{\Delta_{00}+\Delta_{01}}$, $c_H^1(\pi, \Delta) := \frac{1}{\Delta_{00}}$, and $c_C(\pi, \Delta) := \frac{\pi_{01}+\pi_{11}}{\Delta_{01}}$, respectively. The cost of the optimal scheme is equal to the lower envelope of these costs, $\min\{c_E^1, c_M, c_H^1, c_C\}$.*

*One of the following four schemes is optimal in scenario 2 (monitor $i$ receives signal $s_i$; the agent can bribe either one but not both monitors): $S_M$, $S_C$, or*

$$S_E^2 := \left( T_E, \begin{bmatrix} 0 & 0 \\ 0 & \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} \end{bmatrix} \right)$$

$$S_H^2 := \left( T_H, \begin{bmatrix} 0 & 0 \\ \frac{1}{\Delta_{00}} & 0 \end{bmatrix}, \begin{bmatrix} 0 & \frac{1}{\Delta_{00}} \\ 0 & 0 \end{bmatrix} \right)$$

*The costs of the schemes $S_E^2$ and $S_H^2$ are $c_E^2(\pi,\Delta) := \frac{1+\pi_{11}}{\Delta_{00}+\Delta_{01}+\Delta_{10}}$ and $c_H^2(\pi,\Delta) := \frac{1-\pi_{11}}{\Delta_{00}}$ respectively. The cost of the optimal scheme is equal to the lower envelope of the costs, $\min\{c_E^2, c_M, c_H^2, c_C\}$.*

*One of the following four schemes is optimal in scenario 3 (monitors $i$ receives signal $s_i$; the agent can bribe one, other or both of them): $S_E^2, S_M, S_H^1$ or $S_C$. The cost of the optimal scheme is equal to the lower envelope of the costs, $\min\{c_E^2, c_M, c_H^1, c_C\}$.*

*Proof.* Every optimal scheme has a transfer function defined by the set $R$ in lemma 3. The set $R$ can take $4! = 24$ possible values — one for every subset of $S$. For each subset, proposition 1 tells us how to calculate the corresponding wages, so we can easily calculate and compare the costs of the corresponding schemes.

There are four subsets of size 1. The sets $\{(0,0)\}$ and $\{(0,1)\}$ yield the hard and cheeky reward functions $T_H$ and $T_C$ respectively. The set $\{(1,0)\}$ is symmetric to $T_C$, but not as cheap because it rewards the superior evidence, which creates worse incentives than rewarding the inferior evidence (as $T_C$ does). The set $\{(1,1)\}$ is not feasible because the assumption that $\Delta_{11} < 0$ would yield a negative transfer $T(1,1) = 1/\Delta_{11} < 0$, violating limited liability.

There are six subsets of size 2. The set $\{(0,0),(0,1)\}$ yields the moderate reward function $T_M$. The set $\{(0,0),(1,0)\}$ is symmetric but more expensive because it rewards superior evidence (because $\Delta_{01} > \Delta_{10}$ implies $\frac{1}{\Delta_{00}+\Delta_{01}} < \frac{1}{\Delta_{00}+\Delta_{10}}$). The set $\{(1,1),(1,0)\}$ yields a transfer function that costs $\frac{\pi_{11}+\pi_{10}}{\Delta_{11}+\Delta_{10}}$, which is strictly more than the cost of the cheeky scheme because $\Delta_{11} < 0$. So it cannot be optimal. A similar argument applies to $\{(1,1),(0,1)\}$. The 'diagonal' function that rewards in outcomes $(0,1)$ and $(1,0)$ may be feasible, but in all cases it is weakly more expensive than the easy reward scheme. Both are implemented by the same wages, but the easy scheme offers better incentives because it rewards the most exonerating outcome $(0,0)$. The 'diagonal' function that rewards in outcomes $(0,0)$ and $(1,1)$ may be feasible, but in all cases it is strictly more expensive than the hard reward scheme because the latter does not pay rewards in the incriminating outcome $(1,1)$.

There are four subsets of size 3. One of them yields the easy reward function. The set that rewards all outcomes except $(0,0)$ cannot satisfy (IC). The set that rewards all outcomes except $(0,1)$ may be feasible if $\Delta_{01} < 0$, but is always more expensive than the moderate reward scheme. The same is true for the set $\{(0,1),(1,0),(1,1)\}$.

The subsets of size 0 and size 6 entail all transfers being equal, which does not satisfy (IC). Therefore these cannot give solutions.

The costs are easily calculated using proposition 1. It remains to give examples to show that there exist parameters for which each scheme is optimal. The easy scheme is optimal under the parameters $\pi_E := \begin{bmatrix} 0.3 & 0.3 \\ 0.3 & 0.1 \end{bmatrix}$ and $\tau_E := \begin{bmatrix} 0 & 0.1 \\ 0.1 & 0.8 \end{bmatrix}$, which yield the differences These parameters yield costs

$$c_E^1 = \frac{10}{7} < c_E^2 = \frac{11}{7} < c_M = c_C = 2 < c_H^2 = 3 < c_H^1 = \frac{10}{3}.$$

The moderate scheme is optimal under the parameters $\pi_M := \begin{bmatrix} 0.8 & 0 \\ 0.2 & 0 \end{bmatrix}$ and $\tau_M := \begin{bmatrix} 0 & 0.1 \\ 0.1 & 0.8 \end{bmatrix}$, which yield the differences $\Delta_M := \begin{bmatrix} 0.8 & -0.1 \\ 0.1 & -0.8 \end{bmatrix}$. These parameters yield costs

$$c_M = \frac{10}{9} < c_E^1 = c_E^2 = c_H^2 = c_H^1 = \frac{5}{4} < c_C = 2.$$

The hard scheme is optimal under the parameters $\pi_H := \begin{bmatrix} 0.9 & 0 \\ 0 & 0.1 \end{bmatrix}$ and $\tau_H := \begin{bmatrix} 0 & 0.1 \\ 0.1 & 0.8 \end{bmatrix}$, which yield the differences $\Delta_H := \begin{bmatrix} 0.9 & -0.1 \\ -0.1 & -0.7 \end{bmatrix}$. These parameters yield costs

$$c_H^2 = 1 < c_H^1 = \frac{10}{9} < c_M = \frac{10}{8} < c_E^1 = \frac{10}{7} < c_E^2 = \frac{11}{7} < c_C = \infty.$$

The cheeky scheme is optimal under the parameters $\pi_C := \begin{bmatrix} 0.6 & 0.4 \\ 0 & 0 \end{bmatrix}$ and $\tau_C := \begin{bmatrix} 0.1 & 0 \\ 0 & 0.9 \end{bmatrix}$, which yield the differences $\Delta_C := \begin{bmatrix} 0.5 & 0.4 \\ 0.0 & -0.9 \end{bmatrix}$. These parameters yield costs

$$c_C = 1 < c_E^1 = c_E^2 = c_M = \frac{10}{9} < c_H^2 = \frac{9}{5} < c_H^1 = \frac{10}{5}.$$

$\square$

## 1.5  One or two monitors?

**Theorem 1.**

1. *If the agent cannot collude with both monitors simultaneously, then the principal may be strictly better or strictly worse off when each signal is received by a different monitor.*
2. *Otherwise, the principal is always weakly, and sometimes strictly, better off when both pieces of evidence are accessed by the same monitor.*

*Proof.* The key point to note from proposition 2 is that $c_E^1 \leq c_E^2$ and $c_H^2 \leq c_H^3$. All that remains to prove theorem 1 is to find parameters $\pi$ and $\tau$ where these inequalities are strict.

The principal is strictly better off in scenario 2 than in scenario 1 under the parameters $(\pi_H, \tau_H)$, because then $c_H^2 < \min\{c_E^1, c_M, c_H^1, c_C\}$. But the principal is strictly better off in scenario 1 than in scenario 2 under the parameters $(\pi_E, \tau_E)$, because then $c_E^1 < \min\{c_E^2, c_M, c_H^2, c_C\}$. This proves the first statement.

Lemma 1 shows that the principal is more constrained in scenario 3 than in scenario 1, so if the agent can collude with both monitors then she is weakly better off when both pieces of evidence are accessed by the same monitor. The proof of proposition 2 shows that the principal is strictly better off in scenario 1 than in scenario 3 under the parameters $(\pi_E, \tau_E)$ because then $c_E^1 < \min\{c_E^2, c_M, c_H^3, c_C\}$. $\qquad\square$

It follows from the proof of theorem 1 that one monitor is strictly better when $c_E^1 = \frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}} < \min\{c_H^2, c_M, c_C\}$, no matter whether the agent can bribe one or both monitors. This holds only when $\pi_{01} > \tau_{01}$, $\pi_{10} > \tau_{10}$ and $\pi_{11}$ is small, which means that $\pi_{01}\pi_{10} > \tau_{01}\tau_{10}$. Thus one monitor tends to be preferred if the signals are more positively correlated when the agent takes the bad action than when the agent takes the good action, because this implies that $\pi_{11}\pi_{00}-\pi_{01}\pi_{10}$ is small and $\tau_{11}\tau_{00}-\tau_{01}\tau_{10}$ is big. Conversely, two monitors are strictly better than one only if the agent cannot bribe both monitors, and if $c_H^2 = \frac{1-\pi_{11}}{\Delta_{00}} < \min\{c_E^1, c_M, c_C\}$, which tends to be true in the reverse case.

In the special case where both monitors observe exactly the same information, we get that $\pi_{01} = \pi_{10} = \tau_{01} = \tau_{10} = 0$ which implies that $\Delta_{10} = \Delta_{01} = 0$. Hence $c_H^2 < c_E^1 = c_M = c_H^1 < c_H^2 < c_C$, so the principal strictly prefers two monitors if the agent cannot bribe them both; otherwise, she is indifferent.

## 1.6 Robustness

In this section, I show that the main result (theorem 1) continues to hold under various alternative assumptions.

### 1.6.1 Coalitions without the Agent

In the main text, I have considered schemes that are robust to collusion between the agent and one or both monitors. There is also a risk that that the principal might bribe the monitors to conceal evidence. This might seem counter-intuitive, because I have been focussing on the problem of incentivizing the monitors to truthfully report their evidence. But the reason for having the monitors truthfully report their evidence is not that the principal inherently cares about the information, but rather that the principal has to provide credible incentives for the agent to take the right action. By the time the monitor has received the evidence, the agent has already taken his action, so the principal could gain by bribing the monitors to conceal evidence that would otherwise require the principal to reward the agent. If the agent anticipates this then she has no incentive to take the right action.

For example, the cheeky scheme $S_C$ only pays rewards when monitor 2 reports evidence: it rewards the agent in outcome $(0,1)$ and monitor 1 in outcome $(1,1)$. Therefore the principal would be willing to pay monitor 2 any bribe $b < T^C(0,1)$ (any bribe $b < T^C(1,1)$) to conceal their evidence in outcome $(0,1)$ (in outcome $(1,1)$), and monitor 2 would be willing to accept any bribe $b > 0$. So the cheeky scheme is not principal-bribe proof (as opposed to agent-bribe proof). The only way to make the cheeky scheme bribe proof is for the principal to commit to pay the same total quantity rewards in the outcomes $(0,0)$ and $(1,0)$ (when monitor 2 does not report evidence) as in the outcomes $(0,1)$ and $(1,1)$ (when they do). But this undermines the advantage of the cheeky scheme, which was to avoid paying rewards in outcomes $(0,0)$ and $(1,0)$. If the principal is forced to respect principal-bribe proofness then the cost of the scheme increases from $\frac{\pi_{01}+\pi_{11}}{\Delta_{01}}$ to $\frac{1}{\Delta_{01}}$, which is strictly more than the cost of the moderate scheme, $\frac{1}{\Delta_{01}+\Delta_{00}}$.

In general, a scheme is robust to collusion between the principal and a single monitor if they can never increase their joint surplus by concealing evidence:

$$w_i(s) - (T(s) + \sum_{j=1,2} w_j(s)) \geq w_i(m_i, s_{-i}) - (T(m) + \sum_{j=1,2} w_j(m_i, s_{-i})) \quad \forall m_i \leq s_i \text{ and } i = 1,2$$

$$\iff T(m) + w_{-i}(m_i, s_{-i}) \geq T(s) + w_{-i}(s) \quad \forall m_i \leq s_i \text{ and } i = 1,2.$$

These constraints only rule out the cheeky scheme. For instance, if $i = 2$, $s = (0,1)$ and $m_2 = 0$ then the left side equals 0 and the right side equals $\frac{1}{\pi_{01}}$, so the constraint is violated. A scheme is robust to collusion between the principal and both monitors if:

$$\sum_{j=1,2} w_j(s) - (T(s) + \sum_{j=1,2} w_j(s)) \geq \sum_{j=1,2} w_j(m) - (T(m) + \sum_{j=1,2} w_j(m)) \qquad \forall m \leq s$$

$$\iff T(m) \geq T(s) \qquad\qquad\qquad \forall m \leq s.$$

This says that the agent's reward has to increase when the monitors fail to provide hard, incriminating evidence. Like the single monitor constraints, they rule out the cheeky scheme. However, the scenario 3 easy, moderate and hard schemes respect all of these constraints, so there is no change to the main results. The only impact is to expand the regions where the moderate and easy schemes are optimal.

Another possibility is that the monitors might collude with each other, without the agent or the principal. For example, in the scenario 2 easy scheme $S_E^2$, the monitors get a joint payoff of 0 in the outcome $(1,1)$, but they get a collective payoff of $\frac{1}{\Delta_{00}+\Delta_{01}+\Delta_{10}}$ in the outcomes $(0,1)$ and $(1,0)$. Each monitor would be happy to bribe the other monitor to suppress their evidence, and each monitor would be happy to accept the other's bribe. Therefore, this scheme is not robust to monitor-monitor collusion, and can only be made robust by increasing the monitor's collective payoff in $(1,1)$. But doing so yields the scenario 1 easy scheme, $S_E^1$, which undermines the benefits of the second monitor. In general, a scheme is monitor-monitor bribe proof if

$$\sum_{i=1,2} w_i(s) \geq \sum_{i=1,2} w_i(m)$$

for all $m \leq s$. The easy, moderate and hard schemes respect these constraints. Thus, even if the agent cannot bribe both monitors, a single monitor is still weakly better if there is a risk that the two monitors collude with one another.

## 1.6.2 Limited liability

The assumption the players have limited liability rules out the use of punishments, so the principal has to use rewards to provide incentives. This implies that the agent must earn rent from any feasible scheme. In this section, I discuss the implications of replacing the players' limited liability constraints with voluntary participation constraints. Whereas limited liability constraints require that a player's ex post payoff is weakly positive, a voluntary participation constraint instead requires that their ex ante payoff is positive. This means that punishments are permitted, so long as they are balanced by enough rewards to ensure that, on average, the player is indifferent between the scheme and the status quo.

I begin by considering the relaxed problem in which the agent's (LLA) constraint is replaced with a voluntary participation constraint (VPA) requiring that $\sum_{s \in S} \pi_s T(s) \geq 0$. In this case, the principal can reclaim the agent's rent by charging him a fine (negative transfer) in outcomes where she does not reward him. Proposition 2 shows that one of four schemes solves the principal's problem in the (LLA) case. Proposition 3 says that one of the same four schemes can be modified to solve the principal's problem in the (VPA) case, by subtracting the agent's expected reward from his ex post reward. Formally, if $S = (T, w_1, w_2)$ then define $\text{VP}_A(S) :=$ $(T - \mathop{\mathbb{E}}_{s \sim \pi}[T], w_1, w_2)$.

**Proposition 3.** *Suppose the agent's limited liability constraints are replaced with a voluntary participation constraint*

$$\sum_{s \in S} \pi_s T(s) \geq 0. \tag{VPA}$$

1. *One of the following four schemes is optimal in scenario 1 (monitor 1 receives both signals $s_1$ and $s_2$): $VP_A(S_E^1), VP_A(S_M), VP_A(S_H^1)$, or $VP_A(S_C)$. The costs of these schemes are*

$$\tilde{c}_E^1(\pi, \Delta) := \frac{1 - \pi_{00} - \pi_{01} - \pi_{10}}{\Delta_{00} + \Delta_{01} + \Delta_{10}}$$

$$\tilde{c}_M(\pi, \Delta) := \frac{1 - \pi_{00} - \pi_{01}}{\Delta_{00} + \Delta_{01}}$$

$$\tilde{c}_H^1(\pi, \Delta) := \frac{1 - \pi_{00}}{\Delta_{00}}$$

$$\tilde{c}_C(\pi, \Delta) := \frac{\pi_{11}}{\Delta_{01}}$$

*respectively. The cost of the optimal scheme is equal to the lower envelope of these costs, $\min\{\tilde{c}_E^1, \tilde{c}_M, \tilde{c}_H^1, \tilde{c}_C\}$.*

2. *One of the following four schemes is optimal in scenario 2 (monitor $i$ receives signal $s_i$; the agent can bribe either one but not both monitors): $VP_A(S_M), VP_A(S_C), VP_A(S_E^2)$ or $VP_A(S_H^2)$. The costs of schemes $VP_A(S_E^2)$ and $VP_A(S_H^2)$ are*

$$\tilde{c}_E^2 := \frac{1 + \pi_{11} - \pi_{00} - \pi_{01} - \pi_{10}}{\Delta_{00} + \Delta_{01} + \Delta_{10}}$$

$$\tilde{c}_H^2 := \frac{1 - \pi_{11} - \pi_{00}}{\Delta_{00}}.$$

*The cost of the optimal scheme is equal to the lower envelope of the costs, $\min\{\tilde{c}_E^2, \tilde{c}_M, \tilde{c}_H^2, \tilde{c}_C\}$.*

3. *One of the following four schemes is optimal in scenario 3 (monitors $i$ receives signal $s_i$; the agent can bribe one, other or both of them): $VP_A(S_E^2), VP_A(S_M), VP_A(S_H^1)$ or $VP_A(S_C)$. The cost of the optimal scheme is equal to the lower envelope of the costs, $\min\{\tilde{c}_E^2, \tilde{c}_M, \tilde{c}_H^1, \tilde{c}_C\}$.*

*Proof.* I use the fact that the principal's objective depends on the levels of the transfers, whilst the constraints depend on their differences.

First, note that the voluntary participation constraint has to bind; otherwise, the principal could strictly improve her payoff, without violating any constraints, by reducing all of the agent's transfers by a small enough amount. This means that the agent earns no rent in any optimal solution.

Now define a new variable $\tilde{T}(s) := T(s) - T(1,1)$ to be the difference between the agent's reward in outcome $s$ and his reward in outcome $(1,1)$. Note that $\sum_{s \in S} \pi_s \tilde{T}(s) = \sum_{s \in S} \pi_s T(s) - T(1,1)$ because the $\pi_s$ sum to 1; that $\sum_{s \in S} \Delta_s \tilde{T}(s) = \sum_{s \in S} \Delta_s T(s)$ because the $\Delta_s$ sum to 0; and that $\tilde{T}(s) - \tilde{T}(s') = T(s) - T(s')$ for all $s, s' \in S$. Therefore, we can restate the principal's problem as

$$\max_{\tilde{T}, T(1,1)} - \sum_{s \in S} \pi_s(\tilde{T}(s) + W^j(s; \tilde{T}))$$

$$\text{st.} \sum_{s \in S} \Delta_s \tilde{T}(s) \geq 1 \tag{IC}$$

$$\sum_{s \in S} \pi_s \tilde{T}(s) = T(1,1), \tag{VPA}$$

where $W^j$ is defined by proposition 1 (like before), but with $T$ replaced by $\tilde{T}$. Thus the proof of lemma 3 tells us that $\tilde{T}(s) = \frac{\mathcal{I}(s \in R)}{\sum_{s \in R} \Delta_s}$. This implies that $T(1,1) = \frac{\sum_{s \in R} \pi_s}{\sum_{s \in R} \Delta_s}$, so $T^*(s) = \frac{\mathcal{I}(s \in R) - \sum_{s \in R} \pi_s}{\sum_{s \in R} \Delta_s}$.

The proof of proposition 2 then tells us that the schemes listed in the statement are putative optimal schemes. Their costs are easily derived by subtracting off the rent that the agent receives in the corresponding schemes in proposition 2.

$\square$

Thus relaxing the agent's liability constraint does not affect the solution in any important way. Comparing the costs of the schemes in proposition 3 with those in proposition 2 show that the cost of the easy scheme decreases the most when (LLA) is replaced with (VPA). The easy schemes are cheaper in the one monitor scenario than in the two monitor scenarios, so relaxing the (LLA) constraints only reinforces the conclusion of theorem 1 that one monitor is better than two.

If the monitors' liability constraints are both relaxed to voluntary participation constraints, then neither of them earns any rent. In this case, there is no difference between a single monitor and two monitors — their collective surplus is zero in both cases. Instead, I assume that the monitors have some liability limit $L > 0$, where $L$ is small enough to be binding. For $S = (T, w_1, w_2)$ define $\text{VP}_M(S) := (T, w_1 - L, w_2 - L)$.

**Proposition 4.** *Suppose the monitor's limited liability constraints are replaced by weaker liability limits, $w_i(s) \geq -L$, for some exogenous liability limit $L > 0$, and voluntary participation constraints*

$$\sum_{s \in S} \pi_s w_i(s) \geq 0. \tag{VPM}$$

*One of the following four schemes is optimal in scenario 1 (monitor 1 receives both signals $s_1$ and $s_2$): $VP_M(S_E^1)$, $VP_M(S_M)$, $VP_M(S_H^1)$, or $VP_M(S_C)$. The costs of these schemes are*

$$\check{c}_E^1(\pi, \Delta) := \frac{1}{\Delta_{00} + \Delta_{01} + \Delta_{10}} - \min\{L, \frac{\pi_{11}}{\Delta_{00} + \Delta_{01} + \Delta_{10}}\}$$

$$\check{c}_M(\pi, \Delta) := \frac{1}{\Delta_{00} + \Delta_{01}} - \min\{L, \frac{\pi_{11} + \pi_{10}}{\Delta_{00} + \Delta_{01}}\}$$

$$\check{c}_H^1(\pi, \Delta) := \frac{1}{\Delta_{00}} - \min\{L, \frac{\pi_{11} + \pi_{10} + \pi_{01}}{\Delta_{00}}\}$$

$$\check{c}_C(\pi, \Delta) := \frac{\pi_{01} + \pi_{11}}{\Delta_{01}} - \min\{L, \frac{\pi_{11}}{\Delta_{01}}\}$$

*respectively. The cost of the optimal scheme is equal to the lower envelope of these costs, $\min\{\check{c}_E^1, \check{c}_M, \check{c}_H^1, \check{c}_C\}$.*

*One of the following four schemes is optimal in scenario 2 (monitor $i$ receives signal $s_i$; the agent can bribe either one but not both monitors): $VP_M(S_M)$, $VP_M(S_C)$, $VP_M(S_E^2)$ or $VP_M(S_H^2)$. The costs of schemes $VP_M(S_E^2)$ and $VP_M(S_H^2)$ are*

$$\check{c}_E^2 := \frac{1 + \pi_{11}}{\Delta_{00} + \Delta_{01} + \Delta_{10}} - 2\pi_{11}L$$

$$\check{c}_H^2 := \frac{1}{\Delta_{00}} - (\pi_{01} + \pi_{10})L$$

*The cost of the optimal scheme is equal to the lower envelope of the costs, $\min\{\check{c}_E^2, \check{c}_M, \check{c}_H^2, \check{c}_C\}$.*

*One of the following four schemes is optimal in scenario 3 (monitors $i$ receives signal $s_i$; the agent can bribe one, other or both of them): $VP_M(S_E^2)$, $VP_M(S_M)$, $VP_M(S_H^1)$ or $VP_M(S_C)$. The cost of the optimal scheme is equal to the lower envelope of the costs $\min\{\check{c}_E^2, \check{c}_M, \check{c}_H^1, \check{c}_C\}$.*

This result is similar to proposition 3, but now the cost saving for each scheme is proportional to the probability of rewarding the monitors, rather than to the probability of rewarding the agent. The monitors are most likely to be rewarded in the hard scheme (when the agent is least likely to be rewarded), so the hard schemes become relatively more attractive when the monitors have greater liability limits. Since the hard scheme is weakly cheaper in the two monitor scenarios, and strictly so in scenario 2, this suggests that two monitors are relatively more attractive than one when they have relaxed liability constraints.

If monitor $i$ has some positive liability limit $L_i > 0$, then it is always optimal for the principal to reward the agent when monitor $i$ reports exonerating evidence, because she can reclaim the monitor $i$'s wages. Hence the signal $s_i$ will be used even though it may not be very accurate. Finally, if the agent's and at least one of the monitor's liability constraints are relaxed to voluntary participation constraints, then the principal can incentivize the right action for free with a single monitor.

### 1.6.3 More than two monitors

Suppose there is a set $I = \{1, \ldots, n\}$ of signals, each received by a different monitor $i \in I$. Redefine the signal space $S := \{0, 1\}^n$, the distribution of signals when the right action is taken, $\pi \in \Delta\{0, 1\}^n$ and the distribution of signals when the wrong action is taken, $\tau \in \Delta\{0, 1\}^n$. For subset $J \subseteq I$, define $s_J := (s_j)_{j \in J}$, $s_{-J} := (s_i)_{i \notin J}$, and $\pi_J := \sum_{s|s_J=1} \pi_s$. If the agent can bribe all $n$-monitors, then the principal is better off with a single monitor, for the same reason that she is in the two monitor case. But if the agent cannot bribe all the monitors, then the principal may be better off with more monitors. Conjecture 1 considers the case where the agent can bribe a maximum of one monitor.

**Conjecture 1.**

1. *If all pieces of evidence are accessed by monitor 1, so that the principal must respect the agent-monitor bribery constraints*

$$T(s) + w_1(s) \geq T(m) + w_1(m),$$

*for all $m \leq s$, then*
   (a) *any feasible scheme $(T, (w_i)_{i \in I})$ must satisfy*

$$w_1(s) = \max\{T(s) - T(m_J, s_{-J}) \mid m_J \leq s_J\};$$

   (b) *the optimal scheme has a transfer function of the form $T(s) = \mathcal{I}(s \in R)/\sum_{s \in R} \Delta_s$, where $R \subset S$ contains an outcome $s$ only if it contains the lower contour set of $s$;*
   (c) *the cost of the optimal mechanism is $\frac{1}{\sum_{s \in R} \Delta_s}$.*
2. *If each signal is received by a different monitor and the agent can only bribe a single monitor, so that the principal must respect the agent-monitor bribery constraints*

$$T(s) + w_i(s) \geq T(0, s_{-i}) + w_i(0, s_{-i}),$$

*for all $i \in I$, then*
   (a) *any feasible scheme $(T, (w_i)_{i \in I})$ must satisfy*

$$w_i(s) \geq \max\{0, T(s) - T(0, s_{-i})\}$$

*for all $s \in S$;*

*(b) the optimal scheme has a transfer function of the form as in* (1b).

*(c) the cost of the optimal mechanism is* $\frac{\sum_{s \in R} \pi_s}{\sum_{s \in R} \Delta_s}$.

The conjecture says that the optimal mechanism may be easy (if $R = S \setminus \{(1,\ldots,1)\}$), hard (if $R = \{(0,\ldots,0)\}$), or somewhere in between (if, up to reordering, $R$ contains $(1,0,\ldots,0)$ but not $(1,\ldots,1,0)$). Like the two monitor case, each monitor gets rewarded when they alone report evidence, but not when others report evidence. The cost savings from hiring more monitors are potentially large if the optimal set $R$ is small. E.g. if the optimal scheme is a 'hard' one that only rewards the outcome $\mathbf{0} = (0,\ldots,0)$ where all $n$ monitors report exonerating evidence, then the cost of the $n$-monitor scheme is $\pi_{\mathbf{0}} < 1$ times the cost of the one-monitor scheme.

### 1.6.4  Signal structures

I have so far restricted attention to the simplest possible information structure with which to cleanly compare one and two monitor scenarios. If the monitors can receive more than two signals, and can pay to manipulate signals, then the result is largely unchanged — the principal rewards the agent when the monitors report exonerating evidence, and rewards the monitors when they report hard, incriminating evidence. If there is a chance of receiving hard, exonerating evidence and principal-monitor collusion is not possible, then the principal need not reward the monitors at all because they have no conflict of interest. But if principal-monitor collusion is possible then the principal must reward the agent and the monitors when they report hard exonerating evidence, in order to provide credible incentives for the agent.

The case where both monitors observe both signals is isomorphic to the case where they each observe one signal, and the signals are perfectly correlated. The case where monitor 1 observes both signals and monitor 2 observes a single signal is isomorphic to the case where monitor 1 observes a signal with four realisations, $s_1 = (0,0), (0,1), (1,0)$ and $(1,1)$, monitor 2 observes a binary signal $s_2 = 0$ or $1$, and $\mathbb{P}[s_1 \in \{(0,0),(1,0)\}, s_2 = 1] = \mathbb{P}[s_1 \in \{(0,1),(1,1)\}, s_2 = 0] = 0$.

## 1.7  Applications

In this section I discuss how the previous results apply to three different contexts: journalism, whistle-blowing and joint financial auditing.

Journalists (the 'monitors') relate verifiable evidence to the public (the 'principal') about the behaviour of firms, politicians or public institutions (the 'agent'). The public pay journalists for publicising useful information, either through direct purchase of media or through increased advertising revenue. The information relayed by journalists can be used to punish or reward

firms and politicians.[4] However, journalists are susceptible to corruption: they may be bribed to suppress evidence or to create fake stories, or they may be threatened for publicising evidence. Although journalists may be punished for reporting fake news, they are not punished for failing to report real news, so in this respect, they have limited liability.

These stylised facts match the main features of the model. Although I do not attempt to model fake evidence, we can think of 'fake' news as noisy, soft evidence. If there are many monitors, and the agent cannot bribe them all, then proposition 3 and conjecture 1 predict that the public will punish the agent whenever at least one journalist reports incriminating evidence, and journalists will be rewarded only when they alone report evidence. Moreover, the size of each monitor's reward should be proportional to the size of the agent's fine. Reality is not quite so stark as this, but it does seem to be true that a journalist's payoff (including revenue from the story and future earnings from improved career prospects) is increasing in the agent's fine (the size of the story), and decreasing in the number of other journalists who cover the story.

Many regulators offer rewards for hard evidence to whistle-blowers. For instance, in the US the False Claims Act (the 'principal') offers a reward to whistle-blowers ('monitors') who provide evidence that can be used to convict a party (the 'agent') of fraud. The whistle-blower's reward is usually between 15 and 25% of the fine enacted on the convicted party. Similarly, various U.S. acts, including the Clean Water and Clean Air Acts allow any U.S. citizen to file a lawsuit against an offending party. In many cases, the citizens bringing these cases receive a de facto reward for doing so since the amount given as compensation for litigation costs often exceeds the actual litigation. Competition regulation also rewards whistle-blowers for obtaining evidence of breaches of competition law because anti-competitive behaviour is generally harmful to other firms and consumers, so it is in the interests of other firms and consumers to enforce the law. Additionally, some regulators offer explicit financial rewards, for example, the UK Competition and Markets Authority offers rewards of up to a hundred thousand pounds to whistle-blowers who supply evidence of anti-competitive behaviour.

In these situations there coexist a team of 'public inspectors' who obtain low quality evidence with high probability, and a set of 'whistle-blowers' who receive highly accurate evidence with low probability. The baseline model of this paper would suggest that only the whistle-blowers are hired and that they are only rewarded when they obtain incriminating evidence. It seems likely that the public inspectors may have greater liability limits than whistle-blowers, in which case section 1.6.2 predicts that the public inspectors will always be hired and that whistle-blowers will only be hired if the public inspector's liability limit is lower than the size of reward required to satisfy the agent's (IC) constraint. However, the feature that no whistle-blowers are rewarded when they all report incriminating evidence is not explained by heterogeneous monitor liability constraints, rather it is a feature of otherwise optimal mechanisms which do

---

4. For example, facts uncovered by journalists at the Guardian and Channel 4 have been used as evidence to fine the Brexit Campaign for breaking electoral rules.

not respect the agent-multi-monitor-coalition constraints (scenario 2). This seems to be a likely explanation because it is in some contexts difficult to imagine an agent colluding with all the possible whistle-blowers (e.g. local citizens of a polluted lake), so it may not be necessary for the mechanism to respect collusion constraints for certain, large coalitions.

Many countries including France, Germany, and Switzerland, require legal entities to undergo joint financial audits Vanstraelen et al. (2009). Joint financial audits involve two independent accountancy firms working together to produce an audit. Both firms are bear responsibility for the final report, and are liable to be sued for damages arising from inaccurate reports, otherwise they are paid a fixed wage. Legal entities may be fined, sued or otherwise punished in proportion to the amount of incriminating evidence reported. Possible collusion between the auditee and auditors is well known to be a problem since the auditee generally chooses the auditors themselves. Indeed, part of the rationale for joint audits is to reduce the costs of deterring collusion. This outcome is consistent with the predictions of section 1.6.2 when the liability limit for both monitors is strictly positive, but not large enough to rely on a single monitor.

## 1.8 Conclusion

Is it better to receive information from a single monitor, or from two strategically independent monitors? The answer depends on the optimal reward structure, which is, in turn, a function of the (primitive) distribution of evidence. If both types of hard evidence are incriminating enough then the "easy" reward scheme is optimal. But then the agent can receive a reward by bribing either one of the two monitors. The principal can only deter bribes by rewarding both monitors, so she is no better off than when both signals are received by a single monitor. If one signal is more incriminating than the other then the "moderate" scheme is optimal. But this scheme only relies on the superior evidence, so the principal is indifferent between a monitor who receives both types of evidence and a monitor who receives superior evidence only. If neither type of evidence is incriminating enough then the "hard" scheme is optimal. If the agent cannot bribe both monitors, then the principal can save money in the two monitor case by only rewarding each monitor for reporting their evidence when the other fails to do so. This is consistent with real world institutions such as journalism and whistleblowing. But it is not robust to larger collusive coalitions. If the agent can bribe both monitors, then this scheme is vulnerable to larger scale cover-ups because neither monitor stands to be rewarded when they both report evidence. The fact that the agent can act collectively with both monitors as easily as he can with one monitor, means that the principal is weakly better off with a single monitor who receives both signals.

Is it reasonable to assume that the agent can bribe two monitors as easily as a single one? A key assumption in this paper is that bribery negotiations take place under common knowledge of the realised signal $s$ and of the player's respective utility functions. But in reality, there may be a degree of uncertainty about precisely what evidence monitors hold, about the legal implications of this evidence, or about each player's "moral cost" of engaging in bribery. The need to negotiate with differing private evaluations of these factors may not be conducive to the formation of larger side-coalitions. If this is the case, then the principal may prefer information to be received by a larger number of monitors so that she can rely more on their pre-existing asymmetric information to deter bribes (which is free for her), and less on rewards (which are expensive). Another possibility is that the principal can endogenously create asymmetric information with the specific objective of undermining coalition formation (e.g. Ortner & Chassang, 2018). I consider this question in chapter 2.

# Chapter 2

# Lemons by Design

## Abstract

We study a problem in which a firm can bribe an inspector to conceal evidence of illegal pollution. We find that the cheapest way to deter bribes is (i) to secretly select either the firm or the inspector to 'win' a reward whenever evidence is reported; and (ii) to give both the firm and the inspector a secret clue about who will win. If the inspector conceals the evidence, then the winner forgoes their reward — i.e. they 'get a lemon'. The distribution of clues is carefully constructed to engineer the worst possible lemons problem in the market for concealment: player $i$ only enters the market if her clue is strong enough to make her believe that player $j$ is the winner, despite knowing that player $j$ only enters if his clue indicates that player $i$ is the winner. But then higher order reasoning leads neither player to enter the market, no matter what clue they receive. Hence, bribery never takes place in equilibrium. As well as deterring bribes cheaply and robustly, this result demonstrates the full extent of contagious adverse selection in bilateral trades.

## 2.1 Introduction

The threat of corruption is a serious concern in many incentive design problems. For example, firms may bribe politicians to grant them contracts without competitive tender; nations may bribe international inspectors to ignore or destroy evidence of greenhouse gas emissions; and criminals may bribe judges to grant light sentences. Corruption not only reallocates welfare unfairly, but also reduces aggregate welfare by creating perverse incentives. The UN Security Council estimates that corruption in the form of bribes, money laundering and tax evasion, directly reduces world GDP by 5% annually (United Nations, 2018). However this

figure only measures realised corruption, it does not measure the cost of resorting to 'second best' contracts which anticipate the possibility of corruption. In the U.S., approximately 10% of wages are paid to workers whose primary responsibility is monitoring, directly or indirectly, the actions of someone else.[1]

In this paper, we focus on corruption relating to the suppression of evidence by colluding parties. For example Duflo et al. (2013) find evidence that factory owners in India pay pollution inspectors to report that their factories are compliant with regulations, when in fact they are not. Corruption occurs because the factory and the inspector can generate a joint surplus by reporting compliance rather than non-compliance. If corrupt behaviour is detectable and punishable with high enough fines, then the threat of large punishment can deter corruption for free (e.g. Becker (1968); von Negenborn and Pollrich (2020)). But in cases where large punishments are not available, deterring corruption is costly because it requires the government to pay rewards for non-corrupt behaviour. For example, Duflo et al. (2013) incentivise factory inspectors to report truthfully by paying them a reward (efficiency wage) large enough to subsume the joint surplus generated by colluding with the factories. They estimate that it would cost $1300$ per year to enforce compliance for a single small to medium scale factory with high pollution potential.[2] Reducing the cost of these rewards is important both on the intensive margin, because existing regulations that are already being enforced can be enforced at a lower cost; and on the extensive margin, because regulations that are currently too expensive to be enforced, can be made affordable.

Our contribution is to show how information design, together with mechanism or 'transfer' design, can be used to reduce the cost of paying rewards, whilst still implementing compliance. We *engineer a market failure* in the market for bribes by giving the players private information about their transfer so as to create *the worst possible lemons problem*. In this mechanism, each player receives a random number between 0 (lemon) and 1 (peach). The player with the higher number receives a reward (or amnesty, in the case of the firm) for reporting incriminating evidence, so they do not want to conceal evidence. Hence, neither player wants to conceal evidence if they believe that the other player has the lower number (i.e. the lemon). Although it is common knowledge that there is always surplus from corrupt deals, each player is too paranoid to agree to any corrupt deal because higher order reasoning leads them to believe that the other player would only accept a given deal if they received a 0. Hence no concealment takes place.

———

1. Using data from the Occupational Employment and Wage Statistics (OEWS) programme and ONET occupation descriptions.
2. Reported capital investment less than US $ 2 million.

This approach builds on earlier work by Ortner and Chassang (2018) and von Negenborn and Pollrich (2020). They develop endogenous one-sided informational frictions as a means of deterring collusion, but our informational friction is two-sided, and hence more severe. Moreover, our solution accommodates imperfect monitoring (unlike Ortner and Chassang (2018)) and limited liability (unlike von Negenborn and Pollrich (2020)). Our result also provides answers to some open problems identified by Carroll (2016): we find that contagion in higher order beliefs *does* limit the amount of surplus that a pair of agents can obtain when they have transferable utility (in contrast to the case without transferable utility), and we provide an example of a 'worst case' information structure for the transferable utility case. We discuss the relevance of our results to the existing literature in more detail in section 2.6.

The outline of the paper is as follows. Section 2.2 describes the role of the firm and the inspector without intervention from the government. Section 2.3 shows how a regular can use transfers to incentivise the firm to comply with the law, first when bribery is impossible, and then in the case where the firm can bribe the inspector. We introduce the information design in section 2.4 by means of a simple example. Section 2.5.1 formally describes the government's incentive design problem when both transfer design and information design are available as tools. Section 2.5.2 presents our main result: a characterisation of the cheapest scheme that provides adequate incentives for the firm to comply when bribery is possible. The proof of our main results is given in section 2.5.3. Finally, Section 2.6 compares our results with the related literature and section 2.7 concludes with some extensions for future research.

## 2.2 Environment

Finn, the risk neutral firm, chooses whether to pollute, $p$, or comply $c$. Other things equal, polluting yields Finn a payoff of 1, whereas complying yields a payoff of 0. Ina the risk neutral inspector inspects Finn's firm. If Finn chooses to pollute then Ina obtains evidence of pollution with probability $\pi_p$; otherwise, she obtains no evidence. If Finn chooses to comply then she obtains evidence of pollution with probability $\pi_c < \pi_p$. The fact that $\pi_p$ can be strictly less than 1 implies that her monitoring technology is imperfect — she may fail to obtain evidence of pollution even though Finn has been polluting. Similarly, if $\pi_c$ is strictly greater than 0 then she may find evidence that Finn has been polluting, even though he has been compliant. The requirement that evidence is more likely to arise when Finn does pollute than when he complies ensures that evidence is indicative of pollution. If Ina obtains evidence of pollution then she can either report it to the government, or she can keep silent. We assume that evidence cannot be fabricated, so if Ina does not obtain evidence then she has no choice but to keep silent.

The government wants to incentivise Finn to comply, and she does so by fining Finn an amount $f_e > 0$ whenever Ina reports evidence of pollution. If Ina doesn't report evidence then the government fines Finn some amount $f_0$. We refer to the pair $(f_e, f_0)$ as a *scheme*. The fine $f_0$ is necessary because the government doesn't want Finn to shut down completely, so it needs to compensate him for the risk of being fined by mistake when he chooses to comply. More precisely, it must satisfy his voluntary participation constraint,

$$-\pi_c f_e - (1 - \pi_c) f_0 \geq 0,$$

which says that his expected payoff must be weakly positive if he chooses to comply. Since $f_e > 0$, this constraint can only be satisfied by fining Finn a negative amount $f_0 < 0$ whenever Ina reports nothing.

The government must also ensure that Finn prefers to comply than pollute. This is embodied in Finn's incentive compatibility constraint,

$$-\pi_c f_e - (1 - \pi_c) f_0 \geq 1 - \pi_p f_e - (1 - \pi_p) f_0,$$

which says that his expected payoff must be higher if he chooses to comply than if he chooses to pollute. It is more insightful to rewrite this constraint as

$$f_e - f_0 \geq \Pi.$$

where $\Pi := \frac{1}{\pi_p - \pi_c}$. We refer to the difference $f_e - f_0$ as *Finn's incentive to comply* because this is by how much his expected fine falls when he complies. The quantity $\Pi$ is Finn's benefit from polluting divided by the marginal risk of being caught, which we refer to as his risk-adjusted benefit of polluting. Thus Finn's incentive compatibility constraint says that the size of his incentive to comply must exceed his *risk-adjusted* benefit of polluting.

Finally, since Finn's incentive compatibility constraint will ensure he complies in equilibrium, the government's expected cost of a scheme $(f_e, f_0)$ is $-\pi_c f_e - (1 - \pi_c) f_0$. Her objective is to minimise her cost subject to the voluntary participation and incentive compatibility constraints.

## 2.3  Benchmark: transfer design (complete information)

Here we show that incentivising compliance is free when bribes are not possible, but costly when bribes are possible. Thus, the need to deter bribes is the only inefficiency in our model.

The case where the government is unconstrained by the possibility of bribery is the 'first best' case:

**Proposition 5** (First best). *Suppose bribery is not possible. For any constant $k \geq \Pi$, the scheme $(f_e^{FB}, f_0^{FB})$ defined by*

$$f_e^{FB} = (1 - \pi_c)k$$
$$f_0^{FB} = -\pi_c k,$$

1. *satisfies Finn's voluntary participation and incentive compatibility constraints;*
2. *costs $c^{FB} := 0$;*
3. *costs (weakly) less than any other scheme that satisfies Finn's voluntary participation and incentive compatibility constraints.*

*Proof.*

1. Finn's expected payoff is $-\pi_c(1-\pi_c)k - (1-\pi_c)(-\pi_c k) = 0$ so his voluntary participation constraint is satisfied. His incentive is equal to $(1 - \pi_c)k - (-\pi_c k) = k \geq \Pi$ so his incentive compatibility constraint is also satisfied.
2. The cost of the scheme is $c^{FB} := -\pi_c(1 - \pi_c)k - (1 - \pi_c)(-\pi_c k) = 0$.
3. The government's expected cost is exactly equal to Finn's expected payoff and voluntary participation requires Finn's expected payoff to be greater than 0. Therefore no scheme that satisfies voluntary participation can cost less than 0.

$\square$

The schemes described in proposition 5 gives a lower bound on the cost of schemes that incentivise compliance. However, they are not robust to bribes. If incriminating evidence is realised, then bribing Ina to stay silent decreases Finn's fine from $f_e$ to $f_0$. Therefore, Finn will be willing to pay Ina any bribe $b \leq f_e - f_0 = k$. Ina is indifferent between reporting evidence and staying silent, so she will be willing to accept any bribe $b \geq 0$. Thus Ina and Finn can agree to any bribe $0 \leq b \leq k$. If Finn anticipates that he will be able to bribe Ina some amount $b < \Pi$ to stay silent, then polluting and bribing Ina will give him a higher expected payoff than complying, so he will choose to pollute. Consequently, the government needs to deter bribes if it wants to incentivise compliance.

We postpone a detailed description of what it means to be bribery-proof until section 2.5. For the moment, it suffices to say that the government can deter bribes by either reducing the size of Finn's incentive (without violating his incentive constraint), or by paying Ina a reward $r \geq 0$ for reporting evidence. We assume that Ina cannot fabricate evidence, and that she cannot be punished for failing to report evidence (for instance, she may be an employee with limited liability, or a whistle-blower acting of her own volition). Clearly, the government will never want to reward Ina for staying silent, so it is without loss to assume that she pays Ina a reward equal to 0 when she stays silent. Thus a scheme is now defined by a triplet $(f_e, f_0, r)$ and has expected cost $\pi_c(r - f_e) + (1 - \pi_c)(-f_0)$.

Suppose the government continues to impose the first best fines $f_e = (1 - \pi_c)k$ and $f_0 = -\pi_c k$ for some $k \geq \Pi$ and additionally offers Ina a reward $r = k + \epsilon$, for some $\epsilon > 0$. We saw in proposition 5 that these fines satisfy Finn's voluntary participation and incentive compatibility constraints. This scheme also deters bribery because Ina's opportunity cost of staying silent is equal to her reward, so she demands a bribe of at least $k + \epsilon$. But Finn is willing to pay a bribe of at most $k$, so there are no bribes that are mutually agreeable to both Ina and Finn. The cost of this scheme is equal to Ina's expected reward, which is $\pi_c(k + \epsilon)$. This cost is minimised by choosing $k = \Pi$ and $\epsilon$ to be as small as possible. There is no 'smallest' $\epsilon > 0$ so the optimal scheme does not exist. To avoid this technical difficulty, we will be content to say that a scheme deters bribes if it can be made to deter bribes by adding any $\epsilon > 0$ to Ina's reward — we address this point in more detail in section 2.5.1. We refer to the resulting scheme as 'third best' because it is the cheapest scheme among all transfer-only schemes that satisfy Finn's voluntary participation and incentive compatibility constraints; but it is not as cheap as the schemes that utilise information design in sections 2.4 and 2.5.

**Proposition 6** (Third best). *The scheme* $(f_e^{TB}, f_0^{TB}, r^{TB})$ *defined by*

$$f_e^{TB} = (1 - \pi_c)\Pi$$
$$f_0^{TB} = -\pi_c\Pi$$
$$r^{TB} = \Pi,$$

1. *deters bribes and satisfies Finn's voluntary participation and incentive compatibility constraints;*
2. *costs* $c^{TB} := \pi_c\Pi$;
3. *costs less than any other transfer-only scheme that deters bribes and satisfies Finn's voluntary participation and incentive compatibility constraints.*

*Proof.* We have already shown that this scheme satisfies Finn's voluntary participation and incentive compatibility constraints, deters bribes (up to a constant $\epsilon$), and costs $\pi_c\Pi$. It remains to show that no transfer-only scheme costs less. If $f_e - f_0 > r$ then Ina and Finn both strictly benefit by exchanging any bribe $r < b < f_e - f_0$. Therefore any scheme that deters bribes must have $r \geq f_e - f_0$. Any incentive compatible scheme must have $f_e - f_0 \geq \Pi$, therefore any incentive compatible scheme that deters bribes must have $r \geq f_e - f_0 \geq \Pi$. Any scheme that satisfies voluntary participation must have an expected fine of less than 0. Therefore the expected cost of the scheme must be at least $\pi_c r \geq \pi_c\Pi$. □

The intuition is that, without recourse to information design, any bribery proof scheme must destroy all the joint surplus that Ina and Finn generate when Ina stays silent, so we must have $f_e - f_0 - r \leq 0$. But incentive compatibility requires Finn's surplus to be at least $\Pi$, so Ina's reward must be at least $\Pi$.

## 2.4   Illustrative examples: one-sided adverse selection

Paying rewards is necessarily expensive, but information design is free. Akerlof (1970) famously shows how private information in the market for used cars can cause the market to break down. In his model, a buyer is willing to pay a good price for a 'peach' (good car), but he isn't willing to pay anything for a 'lemon' (bad car). A seller is willing to sell a peach for a good price, but he is willing to sell a lemon for any price. The buyer and the seller would be able to generate surplus by trading a peach, but for the fact that only the seller knows whether the car is a peach or a lemon. This creates an adverse selection or 'lemons' problem: the buyer knows that the seller will accept a good price whether she has a peach or a lemon, so the buyer will not want to pay a good price if the proportion of lemons in the market is too high, or if he does not value the peach enough to counteract the risk of buying a lemon. But he won't want to pay anything less than a good price either, because the seller will only accept less than a good price if she has a lemon, which the buyer doesn't value at all. Thus the whole market can collapse, even though there are potential gains from trade.

In this section, we give two examples that illustrate how information design can be used to engineer lemons problems for Ina and Finn. These lemons problems make it more difficult for Ina and Finn to negotiate a bribe, and hence enables the government to deter bribes at a lower cost. Section 2.4.1 presents the simplest scheme with endogenous private information, section 2.4.2 presents the cheapest scheme with endogenous private information for one player only. Section 2.5 presents the cheapest scheme with endogenous private information for both players.

### 2.4.1   The fair coin toss scheme (informed firm)

Before Ina receives any evidence, the government and Finn flip a fair coin with a lemon ($L$) on one side and a peach ($P$) on the other. Finn observes the outcome of the coin toss, whilst Ina does not, so we refer to the outcome as Finn's 'private message'. If the coin comes up lemons then Finn's fines are $f_e(L) = (1 - \pi_c)k$ and $f_0(L) = -\pi_c k$, and Ina's reward is $r(L) = k$, where, as before, $k$ is a constant. In this outcome, *bribery is a lemon for Ina* because she is better off taking her reward of $k$ than any bribe $0 \leq b \leq k$ that Finn is willing to pay. Similar to Akerlof's model, staying silent generates no joint surplus in the lemon state. However, if Finn has a peach message then Finn's fines are reduced to $f_e(P) = \frac{f_e(L)}{2}$ and $f_0(P) = \frac{f_0(L)}{2}$, and Ina's reward is reduced to $r(P) = 0$. In this outcome, *bribery is a peach for Ina* because any bribe $b > 0$ gives her a strictly higher payoff than reporting evidence. Just like Akerlof's model, there is a strictly positive joint surplus of $f_e(P) - f_0(P) = \frac{k}{2} > 0$ in the peach state. Thus Ina's situation is analogous to Akerlof's used car buyer, the only difference being that Akerlof's buyer faces uncertainty about the value of trading, whereas Ina faces uncertainty about the value of *not* trading.

It is easy to verify that the coin toss scheme satisfies Finn's voluntary participation constraint for any choice of $k$. However, the fact that Finn's fine's are reduced in the peach state mean that his expected incentive is now only

$$\frac{1}{2}(f_e(L) - f_0(L)) + \frac{1}{2}(f_e(P) - f_0(P)) = \frac{1}{2}\left[(1 - \pi_c)k + \pi_c k\right] + \frac{1}{2}\left[(1 - \pi_c)\frac{k}{2} - \pi_c\frac{k}{2}\right] = \frac{3}{4}k,$$

so we need to choose $k$ so that $\frac{3}{4}k \geq \Pi$ in order to satisfy this incentive compatibility constraint. The expected cost of this coin toss scheme is

$$\frac{1}{2}\left[\pi_c(r(L) - f_e(L)) + (1 - \pi_c)(-f_0(L))\right] + \frac{1}{2}\left[\pi_c(r(P) - f_e(P)) + (1 - \pi_c)(-f_0(P))\right]$$

$$= \frac{1}{2}\left[\pi_c\{k - (1 - \pi_c)k\} + (1 - \pi_c)\pi_c k\right] + \frac{1}{2}\left[-\pi_c(1 - \pi_c)\frac{k}{2} + (1 - \pi_c)\pi_c\frac{k}{2}\right]$$

$$= \pi_c\frac{k}{2},$$

which is strictly increasing in $k$. Hence the cost is minimised by choosing $k$ to be as small as possible, namely $\frac{4}{3}\Pi$. Doing so yields an expected cost of $\pi_c\frac{2}{3}\Pi$, which is strictly less than the cost of the third best scheme, $c^{\text{FB}} = \pi_c\Pi$.

Despite being cheaper than the third-best scheme, the coin toss scheme still deters all possible bribes (for any choice of $k$) because, just like Akerlof's seller, Finn adversely selects to offer bribes when he has a lemon message, as we now show. First, consider small bribes $b < \frac{k}{2}$. Finn is always willing to pay small bribes because the size of his incentive is weakly greater than $\frac{k}{2}$ in both the peach and the lemon state; specifically, $f_e(L) - f_0(L) = k > f_e(P) - f_0(P) = \frac{k}{2} > b$. Therefore Ina's expected reward is $\frac{1}{2}r(L) + \frac{1}{2}r(P) = \frac{k}{2}$. But this is greater than the size of the bribe, so she will prefer to report the evidence than to remain silent and take a small bribe. Now consider big bribes $\frac{k}{2} < b < k$. Finn pays big bribes if and only if he has a lemon, because the size of his incentive is greater than the bribe in the lemon state, but lower in the peach state. Specifically, $f_0(L) - f_e(L) = k > b > \frac{k}{2} = f_0(P) - f_e(P)$. Therefore, Ina does not accept big bribes because $r(L) = k > b$, so she is better off staying silent in the lemon state. It is clear that Finn will never pay fines strictly greater than $k$ and Ina will never accept bribes strictly less than 0. Thus the only cases left to consider are the knife-edge cases where the bribe exactly equals $k$ or $\frac{k}{2}$, but, no matter what Finn does, Ina is weakly better off rejecting these bribes, so we can make her strictly better off rejecting them by adding some arbitrarily small amount $\epsilon > 0$ to her rewards.

Thus we have proved the following proposition:

**Proposition 7** (Coin toss scheme). *Suppose Finn receives a private message $x_F$ equal to either $L$ (for lemon), or $P$ (for peach). The scheme $(q^{CT}, (f_e^{CT}, f_0^{CT}, r^{CT}))$ defined by*

$$f_e^{CT}(L) = (1 - \pi_c)\frac{4}{3}\Pi \qquad\qquad f_e^{CT}(P) = (1 - \pi_c)\frac{2}{3}\Pi$$

$$f_0^{CT}(L) = -\pi_c\frac{4}{3}\Pi \qquad\qquad f_0^{CT}(P) = -\pi_c\frac{2}{3}\Pi$$

$$r^{CT}(L) = \frac{4}{3}\Pi \qquad\qquad r^{CT}(P) = 0$$

$$q^{CT}(L) = \frac{1}{2} \qquad\qquad q^{CT}(P) = \frac{1}{2},$$

*where $q^{CT}(x_F)$ denotes the probability of the message $x_F$,*

1. *deters bribes and satisfies Finn's voluntary participation and incentive compatibility constraints;*
2. *costs $c^{CT} := \frac{2}{3}\pi_c\Pi$.*

The fair coin toss is the simplest scheme with private information, but not the cheapest. Two variations on the coin toss scheme are possible. Firstly, instead of using a fair coin, the government can use a biased coin which puts a lower weight on the (costly) lemon message. Doing so increases the fraction of peaches in the market at low bribe levels, so the government must further reduce Finn's incentive in the peach state so as to prevent peaches from entering the market at low bribe levels. The optimal biased coin toss scheme costs strictly less than the fair coin toss scheme. Secondly, the government can let Ina observe the outcome of the coin toss instead of Finn. This yields an 'informed inspector' coin toss scheme. We develop this further in the following subsection.

### 2.4.2 The one-sided informed inspector scheme

Coin toss schemes give the informed player a binary signal, which is the minimal amount of private information possible. Here we present a scheme that gives Ina the inspector a whole continuum of possible messages. The reasons for presenting this particular scheme are three-fold. Firstly, it demonstrates that the government can attain the first best outcome if it can use infinitely large fines. We consider the case of infinitely large fines to be unrealistic, so this fact motivates us to consider cases where fines are bounded. Secondly, it is the optimal one-sided scheme (when the bound on fines is not too small[3]), so the fact that our two-sided scheme costs strictly less than it motivates our interest in two-sided schemes. Thirdly, it provides the best basis for comparing our main result with previous literature (Ortner & Chassang, 2018). We expand more on this latter point in section 2.6.

Suppose Ina receives a private message $x_I \in [0, 1]$ with density $q(x_I)$.

---

3. If the bound on fines is small enough then the one-sided informed inspector scheme is undercut by a one-sided informed firm scheme.

**Proposition 8** (Informed Inspector). *For any constant $k \geq \Pi$, the scheme $(q^{II}, (f_e^{II}, f_0^{II}, r^{II}))$ defined by*

$$q^{II} \text{ uniform on } [0,1]$$

$$f_e^{II}(x_I) = \begin{cases} 0 & \text{if } x_I \leq \frac{k-\Pi}{k} \\ (1-\pi_c)k & \text{otherwise} \end{cases}$$

$$f_0^{II}(x_I) = \begin{cases} 0 & \text{if } x_I \leq \frac{k-\Pi}{k} \\ -\pi_c k & \text{otherwise} \end{cases}$$

$$r^{II}(x_I) = \begin{cases} 0 & \text{if } x_I \leq \frac{k-\Pi}{k} \\ k - \frac{k-\Pi}{x_I} & \text{otherwise,} \end{cases}$$

1. *deters bribes and satisfies Finn's voluntary participation and incentive compatibility constraints;*
2. *costs $c^{II} := \pi_c \left[ \Pi + (k-\Pi)\ln\left(1 - \frac{\Pi}{k}\right) \right] \xrightarrow{k \to \infty} 0$;*
3. *costs less than any other one-sided, informed inspector scheme that satisfies Finn's voluntary participation and incentive compatibility constraints.*

*Proof.* We prove here that the informed inspector scheme deters bribes. The rest of the proof is given in appendix A.1.1.

Consider a bribe $b$. Ina will agree to the bribe if and only if she receives a message for which her reward $r(x_I)$ is less than the bribe $b$. We have $r(x_I) = k - \frac{k-\Pi}{x_I}$ so Ina agrees to the bribe if and only if $b \geq k - \frac{k-\Pi}{x_I}$, or equivalently, $x_I \leq \frac{k-\Pi}{k-b}$. This is an example of a cutoff strategy with cutoff equal to $\frac{k-\Pi}{k-b}$. Finn's expected incentive, conditional on Ina's cutoff strategy, is equal to

$$\mathbb{E}\left[ f_e^{II}(x_I) - f_0^{II}(x_I) \,\middle|\, b \geq k - \frac{k-\Pi}{x_I} \right]$$

$$= \mathbb{P}\left[ x_I \leq \frac{k-\Pi}{k} \,\middle|\, x_I \leq \frac{k-\Pi}{k-b} \right] 0 + \mathbb{P}\left[ x_I \geq \frac{k-\Pi}{k} \,\middle|\, x_I \leq \frac{k-\Pi}{k-b} \right] k$$

$$= \min\left\{ 1, \frac{\frac{1}{k-b} - \frac{1}{k}}{\frac{1}{k-b}} \right\} k$$

$$= \min\{k, b\},$$

so he is indifferent about accepting bribes less than $k$, and strictly prefers to reject bribes greater than $k$. If $b = 0$, then Ina strictly prefers to take her reward if her message is $x_I > 1 - \frac{\Pi}{k}$, otherwise she is indifferent. Therefore, Finn's conditional expected incentive is equal to $0$, so he is indifferent as well. In all cases, the government can deter Ina and Finn from exchanging the zero bribe at an arbitrarily small cost (e.g. by adding $\epsilon$ to Ina's reward). $\qquad \square$

The key feature of this scheme is that the distribution of rewards is chosen so that Finn's probability of facing a peach conditional on a given bribe increases in proportion to the size of the bribe, so as to keep him indifferent about accepting the bribe. In other words, 'peach inspectors' enter the market at the highest rate possible without giving Finn a strict preference to enter the market. In the limit, Finn's incentive, $k$, becomes arbitrarily large with vanishing probability, which corresponds to the use of extreme incentives in Becker (1968). By contrast, Ina's reward never exceeds $\Pi$. Since the government only has to pay Ina with the same vanishing probability that she punishes Finn, Ina's expected reward can be made arbitrarily small. This in turn means that the cost of the scheme approaches the first best cost, so no scheme can do better.

However, there is still scope for improvement because transfers are bounded in most practical applications. We show in appendix A.1.2 that, when the $k$ is restricted to be small enough, there is an informed firm scheme that costs less than the optimal informed inspector scheme. In section 2.5 we show that the government can create a two-sided adverse selection problem by jointly designing private information for both the firm and the inspector, and that doing so deters bribes at strictly lower costs than all one-sided schemes, for any bound on transfers.

## 2.5 Main result: two-sided lemons

In this section, we formally model the government's problem with both transfer design and information design. We then describe the players' bribery negotiations and show that the government can restrict attention to bribery-proof schemes. Our main result characterises an optimal bribery-proof scheme.

### 2.5.1 The government's problem

The government commits to a scheme $\mathcal{S} := (q, (f_e, f_0, r))$, where $q$ is a public distribution over private messages $x = (x_I, x_F) \in [0,1]^2$. Finn observes the message $x_F$ and Ina observes the message $x_I$. The objects $f_e, f_0$ and $r$ are message-contingent transfers. We assume that Ina is protected by limited liability, so the government must choose $r(x) \geq 0$. We also assume that the government cannot use extreme incentives for Finn, this means that it must choose $f_e(x)$ and $f_0(x)$ so that $f_e(x) - f_0(x) \leq \kappa$ for some $\kappa \geq \Pi$. If $\kappa < \Pi$ then it will be impossible to provide incentives for Finn to comply; if $\kappa = \infty$ then the government can achieve the first best with the one-sided scheme described in section 2.4.2. Like before, the government's main objective is to incentivise Finn to choose to comply instead of polluting. This requires it to satisfy an incentive compatibility constraint. It also needs to incentivise Finn to stay in business, which requires it to satisfy a voluntary participation constraint. We say that the scheme $\mathcal{S}$ is *feasible* if it satisfies these three constraints: limited liability, incentive compatibility, and voluntary participation. The government wants to find the cheapest feasible scheme.

After the government has sent Ina and Finn their private messages, Ina may or may not obtain evidence. If she does obtain evidence then Ina is able to commit to stay silent in return for a bribe from Finn. For each possible bribe $b$ we define a 'bribery game' in which Ina and Finn choose to accept the bribe with respective probabilities $\sigma_I^b(x_I)$ and $\sigma_F^b(x_F)$. If they both accept the bribe then Finn pays the bribe $b$ to Ina and Ina stays silent: Ina's payoff is then $b$ and Finn's payoff is $f_0(x) - b$. If either of them rejects the bribe then Ina reports the evidence and gets paid $r(x)$; Finn receives the fine $f_e(x)$. We refer to the bribe-strategy pair $(b, (\sigma_F^b, \sigma_I^b))$ as a *bribe contract*, and we denote it by $\mathcal{C} := (b, (\sigma_F^b, \sigma_I^b))$. We refer to the bribe contract in which Ina always reports evidence and Finn never pays bribes as the *null bribe contract*, denoted $\mathcal{C}_0$. If Ina doesn't obtain evidence then she gets 0 and Finn gets $f_0(x)$.

**Definition 1.** A scheme $\mathcal{S}$ is *bribery proof* if Ina always reports evidence in every Bayes-Nash equilibrium of every bribery game. A scheme $\mathcal{S}$ is $\epsilon$-*bribery proof* if the scheme $(q, (r + \epsilon, f_e, f_0))$ is bribery proof for all $\epsilon > 0$.

In other words, a scheme is $\epsilon$-bribery proof if it can be made bribery proof at an arbitrarily small cost. The reason for introducing $\epsilon$-bribery proofness is purely technical — it ensures that an optimal scheme exists. Without it, we would have to add $\epsilon$ to our schemes to destroy any unwanted equilibria. Another potential solution to this problem would be to do partial implementation, but this would give a trivial solution in our model because the null bribe contract is always an equilibrium: rejecting the bribe is always a weak best response if the other player always rejects it.

A natural question arises about how Ina and Finn choose between different possible bribe contracts, but this turns out not to matter because lemma 4 tells us that the government can restrict attention to $\epsilon$-bribery proof schemes.

**Lemma 4.** *For every feasible scheme, there exists an $\epsilon$-bribery proof feasible scheme with the same cost.*

*Proof.* Let $\mathcal{S}$ be a feasible scheme. If $\mathcal{S}$ is not $\epsilon$-bribery proof then there must exist a non-null, bribery equilibrium. A bribery equilibrium is payoff equivalent to a special case of an incentive compatible collusive side contract.[4] Therefore an interim efficient, non-null side contract, $\mathcal{C}$, must exist. The government can create a new scheme $\mathcal{S}'$ that replicates the ex post payoffs of the bribe contract $\mathcal{C}$ under the original scheme $\mathcal{S}$. This new scheme $\mathcal{S}'$ must be $\epsilon$-bribery proof: otherwise, there would exist a non-null, interim efficient side contract $\mathcal{C}'$ that gives either Ina or Finn must a strictly higher payoff under $\mathcal{S}'$, than does the null contract $\mathcal{C}_0$.

---

4. See appendix A.1.3.

To see why such a $\mathcal{C}'$ cannot exist, note that the scheme, contract pair $(\mathcal{S}',\mathcal{C}_0)$ is ex post payoff equivalent to $(\mathcal{S},\mathcal{C})$ so $(\mathcal{S}',\mathcal{C}')$ would interim dominate $(\mathcal{S},\mathcal{C})$. But then we could construct another side contract $\mathcal{C}''$ that is ex post payoff equivalent under $\mathcal{S}$ to $(\mathcal{S}',\mathcal{C}')$. Since $\mathcal{C}'$ is incentive compatible, $\mathcal{C}''$ must also be incentive compatible, contradicting the fact that $(\mathcal{S},\mathcal{C})$ was assumed to be interim efficient. Therefore $\mathcal{S}'$ must be $\epsilon$-bribery proof. $\qquad\square$

The fact that the government can restrict attention to $\epsilon$-bribery proof schemes simplifies the problem dramatically, because it means that bribes need not feature in Ina and Finn's payoffs. The cost of this simplification is that the scheme must satisfy an $\epsilon$-bribery proofness constraint. Formally, the problem is

$$\min_{\mathcal{S}} \mathbb{E}[\pi_c(r(x) - f_e(x)) + (1 - \pi_c)f_0(x)]$$

$$\text{s.t. } \mathbb{E}[\pi_c f_e(x) + (1 - \pi_c)f_0(x)] \geq \mathbb{E}[\pi_d f_e(x) + (1 - \pi_d)f_0(x)] \qquad \text{(IC)}$$

$$\mathbb{E}[\pi_c f_e(x) + (1 - \pi_c)f_0(x)] \geq 0 \qquad \text{(VP)}$$

$$r(x) \geq 0 \text{ and } f_e(x) - f_0(x) \leq \kappa \qquad \text{(LL)}$$

$$\mathcal{S} \text{ is } \epsilon\text{-bribery proof.} \qquad \text{(BP)}$$

We are now ready to present our main result.

### 2.5.2 An optimal scheme

Our main result characterises an optimal scheme with two-sided information design. This scheme endogenously creates a two-sided lemons problem for Ina and Finn. In the informed inspector scheme, Ina's willingness to accept bribes was decreasing in her own message, so she only wanted to agree to bribes when her message was below a certain cutoff. Finn's willingness to accept bribes was increasing in Ina's message, so Ina's choice of cutoff strategy made him unwilling to accept any bribes. In the two-sided scheme, Ina is both informed (about her own message) and uninformed (about Finn's message) so she inherits both of these features. Her willingness to accept bribes is decreasing in her own message, so she continues to adopt a cutoff strategy. But her willingness to accept bribes is increasing in Finn's message, so her optimal choice of cutoff is decreasing in her belief about Finn's message. Her belief about Finn's message depends on his strategy. The distribution of transfers ensures that his best response is also a cutoff strategy, and his optimal cutoff is decreasing in his belief about Ina's message. The distribution of transfers is chosen so that each player wants to choose a cutoff that is slightly below the cutoff that the other player chooses. Therefore there can be no equilibrium in which they both use a positive cutoff, which means that they do not agree to bribes in any equilibrium. Thus, unlike the informed inspector scheme, the two-sided scheme uses contagion in higher order beliefs to amplify the adverse selection problem. This stands in contrast to an earlier result of Carroll (2016), which we describe in more detail in section 2.6.

The functional form of our optimal two-sided scheme is comparatively simple: messages are independently and uniformly distributed, and, in the special case where $\kappa = 1$ and $\Pi = \frac{3}{4}$, transfers are a linear function of the ratio of the two messages, truncated above at 1:

$$r^*(x) = 1 - \min\left\{1, \frac{x_F}{x_I}\right\}$$

$$f_e^*(x) - f_0^*(x) = \min\left\{1, \frac{x_I}{x_F}\right\}.$$

Figure 2.1 shows these transfers in the context of the message space. When Ina's message is lower than Finn's, her reward is equal to 0, so she is eager to accept any bribe. However, Finn gets a reduced incentive of $\frac{x_I}{x_F} < 1$, so staying silent becomes less valuable to him and therefore he is not willing to pay such high bribes. Similarly, When Ina's message is higher than Finn's, Finn's incentive is equal to 1, so he is eager to accept any bribe, but Ina gets a reduced reward of $1 - \frac{x_F}{x_I} < 1$, so she is relatively less willing to accept bribes. The fact that messages are independent in this optimal scheme is not surprising: Cremer and McLean (1988) tell us that a third party would be able extract Ina and Finn's private information for free if their messages were correlated.[5] Having extracted their private information, this third party could then choose a bribe that would be mutually agreeable to Ina and Finn.



**Figure 2.1:** Transfers as a function of messages in the special case $\kappa = 1$, $\Pi = \frac{3}{4}$.

In the general case where the bound on Finn's incentive takes any value large enough to satisfy incentive compatibility (i.e. $\kappa \geq \Pi$), the two-sided scheme takes a similar form, described in theorem 2:

———

5. See appendix A.1.3.

**Theorem 2** (Two-sided adverse selection). *The scheme $\mathcal{S}^* := (q^*, (r^*, f_e^*, f_0^*))$ defined by*

$$q^* \text{ independent and uniform on } [0,1]^2$$

$$r^*(x) = \kappa \left(1 - \min\left\{1, x_F^\lambda/x_I\right\}\right)$$

$$f_e^*(x) = \kappa(1 - \pi_c) \min\left\{1, x_I/x_F^{1/\lambda}\right\}$$

$$f_0^*(x) = -\kappa\pi_c \min\left\{1, x_I/x_F^{1/\lambda}\right\},$$

*where $\lambda = \sqrt{\frac{\kappa}{\kappa - \Pi}} - 1$,*

1. *is feasible;*
2. *costs $c^* := \pi_c \left(\sqrt{\kappa} - \sqrt{\kappa - \Pi}\right)^2$;*
3. *costs less than any scheme that can be approximated by feasible schemes with a finite number of messages.*

The two-sided scheme is a substantial improvement on the optimal one-sided scheme. Table 2.1 compares the cost of the two-sided scheme to the one-sided (informed inspector) scheme (and others) for a range of parameter values. We show in a appendix A.1.4 that the two-sided scheme costs strictly less than the informed inspector scheme at all parameter values and that the cost of the two-sided scheme converges to *half* the cost of the optimal one-sided scheme as $\kappa$ gets large.

**Table 2.1:** The costs of selected schemes for parameters $\pi_p = \frac{2}{3}, \pi_c = \frac{1}{3}, \Pi = 3$.

| Scheme | $\kappa = 3$ | $\kappa = 4$ | $\kappa = 6$ | $\kappa = 30$ | $\kappa = \infty$ |
|---|---|---|---|---|---|
| third best | 1 | 1 | 1 | 1 | 1 |
| fair coin toss (informed firm) | — | 0.667 | 0.667 | 0.667 | 0.667 |
| biased coin toss (informed firm) | 1 | 0.667 | 0.586 | 0.513 | 0.5 |
| biased coin toss (informed inspector) | 1 | 0.75 | 0.5 | 0.1 | 0 |
| one-sided (informed inspector) | — | 0.538 | 0.307 | 0.052 | 0 |
| two-sided | 1 | 0.333 | 0.172 | 0.026 | 0 |

### 2.5.3   Proof of the main result

We prove each of the numbered points in theorem 2 in turn.

**The two-sided scheme is feasible**

The proof has three steps: First, we show that cutoff strategies are optimal; then we calculate each player's expected returns to bribes when the other player plays a cutoff strategy; finally, we show that no bribe is mutually agreeable at any pair of cutoffs. An alternative proof replaces the final step with a demonstration that the best response to any cutoff is a proportionally lower cutoff, with the result that no pair of strictly positive cutoffs can form an equilibrium.

Fix a bribe $b$. The same proof applies to all bribes.

*Step 1.* Ina's opportunity cost of agreeing to the bribe when she receives message $x_I$ is equal to her expected reward conditional on Finn's strategy $\sigma_F$:

$$\mathbb{E}_{x_F}[\sigma_F(x_F)r^*(x)] = \mathbb{E}_{x_F}\left[\sigma_F(x_F)\kappa\left(1 - \min\{1, x_F^\lambda x_I^{-1}\}\right)\right]. \tag{2.1}$$

This expression is continuous and strictly increasing in her own message $x_I$ whenever $\sigma_F(x_F) > 0$ for some $x_F$ (otherwise it equals 0). Therefore, Ina's best response to any $\sigma_F$ is to agree to a bribe $b$ iff her message is below some cutoff $y_I \in [0,1]$. If equation (2.1) is greater than $b$ for all $x_I$ then Ina never wants to agree to the bribe, so she chooses a cutoff of $y_I = 0$. If equation (2.1) is less than $b$ for all $x_I$ then Ina always wants to agree to the bribe, so she chooses a cutoff of $y_I = 1$. In intermediate cases, Ina chooses her cutoff to equal the unique message $x_I$ for which (2.1) is exactly equal to $b$.

The same argument applies to Finn: his expected incentive conditional on Ina's strategy $\sigma_I$ is

$$\mathbb{E}_{x_I}[\sigma_I(x_I)(f_e^*(x) - f_0^*(x))] = \mathbb{E}_{x_I}\left[\sigma_I(x_I)\kappa\min\{1, x_I/x_F^{1/\lambda}\}\right],$$

which is also continuous and strictly decreasing in his own message $x_F$ whenever $\sigma_I(x_I) > 0$ for some $x_I$. Therefore, Finn's best response to any $\sigma_I$ is to agree iff his message is below some cutoff $y_F \in [0,1]$.

*Step 2.* When Ina receives the message $x_I$ and Finn uses cutoff $y_F$, Ina only agrees to the bribe if it is above her expected opportunity cost

$$\mathbb{E}[r^*(x)|x_I, x_F \leq y_F] = \begin{cases} \left(1 - \frac{1}{\lambda+1}\frac{y_F^\lambda}{x_I}\right)\kappa & \text{if } x_I \geq y_F^\lambda, \\ \frac{\lambda}{\lambda+1}\frac{x_I^{1/\lambda}}{y_F}\kappa & \text{if } x_I < y_F^\lambda. \end{cases} \tag{2.2}$$

When Finn receives the message $x_F$ and Ina uses cutoff $y_I$, Finn only agrees to the bribe if it is below his expected opportunity cost

$$\mathbb{E}[f_e^*(x) - f_0^*(x)|x_F, x_I \leq y_I] = \begin{cases} \frac{\lambda}{\lambda+1}\frac{y_I^{1/\lambda}}{x_F}\kappa & \text{if } x_F^\lambda > y_I, \\ \left(1 - \frac{1}{\lambda+1}\frac{x_F^\lambda}{y_I}\right)\kappa & \text{if } x_F^\lambda \leq y_I. \end{cases} \tag{2.3}$$

The full derivation of equations (2.2) and (2.3) are shown in appendix A.1.5.

*Step 3.* To be mutually accepted at the cutuoff outcome $y = (y_F, y_I)$, the bribe $b$ must be greater than Ina's conditional expected reward and below Finn's conditional expected incentive:

$$\mathbb{E}[r^*(x)|x_I = y_I, x_F \leq y_F] \leq b \leq \mathbb{E}[f_e^*(x) - f_0^*(x)|x_F = y_F, x_I \leq y_I].$$

But if $y_F^\lambda \leq y_I$, then

$$\mathbb{E}[r^*(x)|x_I = y_I, x_F \leq y_F] = \frac{\lambda}{\lambda + 1} \frac{y_I^{1/\lambda}}{y_F} \kappa = \mathbb{E}[f_e^*(x) - f_0^*(x)|x_I = y_I, x_F \leq y_F].$$

And if $y_F^\lambda > y_I$, then

$$\mathbb{E}[r^*(x)|x_I = y_I, x_F \leq y_F] = \left(1 - \frac{1}{\lambda + 1} \frac{y_F^\lambda}{y_I}\right) \kappa = \mathbb{E}[f_e^*(x) - f_0^*(x)|x_I = y_I, x_F \leq y_F].$$

In both cases, Ina and Finn are both indifferent about agreeing to the bribe when they receive their cutoff messages $y_I$ and $y_F$. This means that adding any $\epsilon > 0$ to Ina's reward will make her strictly prefer to reject the bribe at her cutoff and for nearby messages. Therefore, the cutoffs $y_I$ and $y_F$ cannot form an equilibrium. Since this argument applies for any pair of positive cutoffs, the bribe will not be accepted in any equilibrium.

*Incentive compatibility* Evaluating equation (2.3) at $y_I = 1$ gives $\mathbb{E}[f_e^*(x) - f_0^*(x)|x_F] = \mathbb{E}[f_e^*(x) - f_0^*(x)|x_F, x_I \leq 1] = \left(1 - \frac{1}{\lambda+1} x_F^\lambda\right)\kappa$. Integrating over $x_F$ gives Finn's ex ante incentive:

$$\mathbb{E}[f_e^*(x) - f_0^*(x)] = \int_0^1 \mathbb{E}[f_e^*(x) - f_0^*(x)|x_F] \, dx_F = \kappa \int_0^1 1 - \frac{1}{\lambda + 1} x_F^\lambda \, dx_F = \kappa \left(1 - \frac{1}{(\lambda + 1)^2}\right).$$

Substituting in $\lambda = \sqrt{\frac{\kappa}{\kappa - \Pi}} - 1$ gives $\mathbb{E}[f_e^*(x) - f_0^*(x)] = \Pi$, as required.

*Voluntary participation*

$$\mathbb{E}[\pi_c f_e^*(x) + (1 - \pi_c) f_0^*(x)] = \quad \mathbb{E}[\pi_c \kappa (1 - \pi_c) \min\{1, x_I/x_F^{1/\lambda}\} - (1 - \pi_c)\kappa \pi_c \min\{1, x_I/x_F^{1/\lambda}\}]$$
$$= 0$$

**The two-sided scheme costs $\pi_c(\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$.**

Evaluating equation (2.2) at $y_F = 1$ gives $\mathbb{E}[r^*(x)|x_I] = \mathbb{E}[r^*(x)|x_I, x_F \leq 1] = \frac{\lambda}{\lambda+1}x_I^{1/\lambda}\kappa$. Integrating over $x_I$ gives Ina's expected reward:

$$\mathbb{E}[r^*(x)] = \int_0^1 \frac{\lambda}{\lambda+1}x_I^{1/\lambda}\kappa\, dx_I$$
$$= \left(\frac{\lambda}{\lambda+1}\right)^2\kappa.$$

Substituting in $\lambda = \sqrt{\frac{\kappa}{\kappa - \Pi}} - 1$ gives $\mathbb{E}[r^*(x)] = (\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$, so the cost of the scheme is is $\pi_c(\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$.

**Every feasible finite scheme costs at least $\pi_c(\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$.**

The government's problem involves deterring all bribes. We show that the cost of deterring the 'worst case' bribe $b^* = \sqrt{\kappa}\left(\sqrt{\kappa} - \sqrt{\kappa - \Pi}\right)$, is equal to the cost of the two-sided scheme. Since deterring the specific bribe $b^*$ is easier than deterring all bribes simultaneously, the cost of deterring $b^*$ gives a lower bound on the cost of deterring all bribes. Therefore no feasible scheme can cost strictly less than the two-sided scheme.

It remains to show that any scheme that deters the bribe $b^*$ costs at least $\pi_c(\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$. This problem is dramatically simplified by a result in Carroll (2016) which shows that we can restrict attention to public schemes. This result is stated in lemma 5.

**Lemma 5** (Corollary of Carroll (2016))**.** *Suppose a finite scheme $\mathcal{S}$ deters a bribe $b$. Then there exists a public, finite scheme, $\mathcal{S}^p$ (i.e. one in which $x_I = x_F$ with probability 1) which deters the bribe $b$ and costs the same as $\mathcal{S}$.*

*Proof.* We first translate our problem into the model of Carroll (2016). The result is then a corollary of his propositions 3.1 and 3.2. The details are given in appendix A.1.6. □

A public scheme deters bribe $b^*$ if and only if, for all $x \in \text{supp}(q)$, either $r(x) > b^*$ or $f_e(x) - f_0(x) < b^*$. An optimal public scheme with incentive equal to $\Pi$ has $r(x) = b^*$ and $f_e(x) - f_0(x) = \kappa$ with probability $\frac{\Pi - b^*}{\kappa - b^*}$; and $r(x) = 0$ and $f_e(x) - f_0(x) = b^*$ with probability $1 - \frac{\Pi - b^*}{\kappa - b^*}$. This scheme has expected reward $\frac{\Pi - b^*}{k - b^*}b^* = (\sqrt{\kappa} - \sqrt{\kappa - \Pi})^2$, hence its cost is equal to $c^*$. Therefore, the two-sided scheme costs less than every finite feasible scheme. It follows from continuity that it costs less than every scheme that can be approximated by a sequence of finite feasible schemes.

## 2.6 Literature

We contribute to an extensive literature on corruption (Baliga & Sjöström, 1998; Laffont & Martimort, 1997; Strausz, 1997; Tirole, 1986). The closest paper to ours is Ortner and Chassang (2018). They are the first (to the best of our knowledge) to study the use of endogenous asymmetric information to deter bribes. They show that a principal (the government) can benefit from paying the monitor (Ina) a random wage (privately observed by the monitor) according to a public distribution, known to the agent (Finn). Doing so endows the monitor with private information about their outside option, and thereby creates an informational-friction in subsequent collusive negotiations between the monitor and a would-be criminal agent. In their model, the agent chooses between being criminal and bribing the monitor on the one hand, or being innocent on the other. Paying the monitor a random wage creates a trade-off for the agent: he can either offer a high bribe which guarantees a high probability of successfully corrupting the monitor, or he can offer a low bribe which guarantees a low probability of successfully corrupting the monitor. The principal saves money by paying random wages because low wage monitors can mimic the high wage monitors and demand high bribes.

Our model differs from theirs in two important respects. Firstly, Ortner and Chassang (2018) assume perfect monitoring so they can rule out bribes on the equilibrium path (when no incriminating evidence arises) without needing to rule them out off the equilibrium path (when the monitor receives incriminating evidence with certainty). By contrast, we allow for monitoring mistakes so we have to consider the impact of bribes both on and off the equilibrium path. If, like Ortner and Chassang, we pay rewards according to a distribution that pays rewards that are smaller than the agent's punishment, then we inevitably get on-path bribery because the agent is always weakly better off accepting bribes smaller than the punishment, and there will be a strictly positive probability that the monitor is willing to accept such bribes. This difficulty motivates our second main departure from their model, which is to endogenise the agent's fines. Doing so allows us to replicate the agent's trade-off in their model, because we can use the changes in the agent's fine to imitate his choice to commit crime or not. This gives our result a qualitatively different interpretation from theirs: our agent faces a lemons problem because his fine depends on the monitor's private information. Despite these differences, our informed inspector scheme (proposition 8), which is the closest to theirs conceptually, produces the same distribution of rewards and has the same cost.[6] We showed in section 2.5.2 that the two-sided scheme costs strictly less than the informed inspector scheme, and costs half as much in the limit as the size of the maximal punishment increases.

_____

6. Garrett, Georgiadis, Smolin, and Szentes (2020) obtain the same distribution as a solution to a similar problem in which an agent chooses their distribution of costs to maximise their information rent.

Another closely related paper to ours is von Negenborn and Pollrich (2020). They also find that engineering a lemons problem is an optimal solution to a mechanism design problem. Our main contribution relative to theirs is that we impose bounds on all transfers, whereas their proposed mechanisms attain the first best by using large rewards and/or punishments. Therefore, they do not need to engineer an optimal lemons problem – any lemons problem would suffice.

Our results also speak to the literature on the robustness of equilibria to contagion.[7] Carroll (2016) obtains an upper bound on the amount of surplus lost due to contagion in a game with two agents either accepting or rejecting a proposed deal, where both agents have private information about the payoff outcomes of the deal. Surprisingly, Carroll finds that contagion does not prevent the agents from realising joint surplus, so long as they have common knowledge that their ex-post joint surplus from the deal is weakly positive. He concludes by asking "What change[s] if we consider ... mechanism[s] that determine not only whether a deal takes place but which deal is chosen? ... Is it possible to describe the worst-case information structure?" (pp. 355–356). Our two-sided scheme entails common knowledge that the ex-post joint surplus from bribery is weakly positive, and yet we find that contagion does play an important role in this scheme. We conclude that contagion does become problematic when the players are trying to negotiate the terms of a deal, because the players' types adversely select which terms to accept. The two-sided scheme has a worst case information structure which leads to all deals being rejected for a particular distribution of payoffs (payoffs are endogenous in our setting, but exogenous in his). This worst case information structure has independent and uniform signals that quantify the severity of the lemons problem faced by the recipient.

Our problem fits into a larger class of general mechanism design problems in which the designer chooses both transfers and information.[8] A particularly relevant and recent paper is Halac, Lipnowski, and Rappoport (2021)'s *Ranking Uncertainty in Organisations*. They show how 'ranking schemes' can create strategic uncertainty and thereby induce a team of workers to exert complementary effort on a project. Ranking schemes are superficially similar to ours in two respects. Firstly, all the players receive a private message. Secondly, the distribution of payoffs is chosen so that work is a dominant strategy for the players with the highest possible message realisation, and each player finds it optimal to work conditional on the belief that all players with the same message or higher will work. Thus, like ours, their scheme produces an inductive chain that causes working to be a higher order best response for all other workers. However, the mechanism underlying their ranking scheme is qualitatively different from ours. Endogenous private information benefits their designer because each worker's incentives to work are strictly concave in their belief that other workers will work. Therefore, a given incentive is created more cheaply by randomising over beliefs. There is no

---

7. See e.g. Kajii and Morris (1997); Morris and Ui (2005).
8. See Bergemann and Morris (2019); Mathevet, Perego, and Taneva (2020); Taneva (2019).

lemons problem in their scheme because workers have complete information about their own payoffs — other workers' types only affect them indirectly through the other workers' decisions to exert effort. By contrast, asymmetric information only benefits us because it inhibits our players from negotiating bribes (which are not considered in Halac et al. (2021)). We engineer a lemons problem by designing a scheme in which each player's payoffs depend directly on the message received by the other player.

## 2.7 Conclusion

We show how information design can be used together with transfer design to deter bribes by engineering a lemons problem. The optimal scheme characterised in our main result, theorem 2, accommodates monitoring errors, costs strictly less than other schemes in the literature, and is relatively simple to implement. This scheme also gives insights into 'worst case' information structures and gives an upper bound on the amount of surplus lost to contagion in bargaining games.

We have shown that the two-sided scheme deters bribes, but bribery contracts are only a special case of an incentive compatible side-contract. Side contracts additionally allow for the possibility of message dependent bribes, $b(x)$, and correlated agreement strategies that depend on both messages, $\sigma(x)$. We conjecture that the two-sided scheme deters all side contracts. One possibility for future research is to prove this by showing that the core of the cooperative game with incomplete information (Forge & Serrano, 2013; Myerson, 2007) induced by the scheme is empty.

Studying the core of the cooperative game induced by the two-sided scheme would also be valuable for extending our results to more than two players. We see this as a particularly promising avenue for future research because it could help us to utilise the information held by potential whistleblowers. Specifically, if the lemons problem can be made disproportionately worse by spreading information across an even larger number of 'inspectors', then we expect to find that the costs of implementing compliance can be further reduced by offering stochastic rewards to whistleblowers. This stands in contrast to the case without information design, where hiring multiple monitors does not help to deter bribery (Stapenhurst, 2019).

Finally, the use of endogenous lemons to deter collusion may have applications beyond monitoring. For instance, it can be used to deter illegal trades, such as weapons, drugs, and human tracking. We also speculate that it could also be used to break up cartels or to deter sub-coalitions of would-be signatories from undermining international agreements.

# Inferring Inequality: Testing for Median-Preserving Spreads in Ordinal Data

## 3.1 Introduction

Health and wellbeing, educational qualifications, standards of sanitation, credit ratings, and perceived corruption are all examples of ordinal variables studied by social scientists. These variables take values that can be ordered, but not quantified. For example, the EUROSTAT Survey on Incomes and Living Conditions asks respondents to rate their happiness on a five-point scale from (1) very bad, to (5) very good. Movements up the scale correspond to an improvement in happiness, but the size of the improvement in moving from, say, 'very bad' to 'bad' may not be the same as that in moving from 'good' to 'very good'. We know that some standard summary statistics, such as the mean and variance, cannot be meaningfully applied to such variables (Stevens, 1946). For instance, Allison and Foster (2004) show that the choice of scale can determine which of two ordered multinomial distributions has the higher variance.

Several authors have responded to this problem by developing purpose-made inequality indices for ordinal variables which are not sensitive to arbitrary scale choices (see Silber & Yalonetzky, 2021, for a review). Many of these indices (Abul Naga & Yalcin, 2008; Apouey, 2007; S. R. Chakravarty & Maharaj, 2015; Kobus & Milos, 2012; Lazar & Silber, 2013; Reardon, 2009) respect the "median-preserving spread" partial ordering of Allison and Foster (2004).[1] A distribution $G$ is a *median-preserving spread* (MPS) of a distribution $F$ ('$F$ and $G$ are ordered') if $F$ and $G$ share a common median and if the probability mass of $G$ lies further away from the median category than $F$'s (i.e. $G$ has 'thicker tails' than $F$). Therefore, $G$ is

---

1. The MPS partial ordering is itself a special case of Mendelson (1987)'s "quantile-preserving spread"; see section 3.6.

deemed more unequal than $F$ according to *all* of these inequality indices (Kobus, 2015), if (and only if) $G$ is an MPS of $F$. This is analogous to the result that two cardinal distributions being ordered by *mean*-preserving spread implies that one is in every sense riskier than the other (Rothschild & Stiglitz, 1970).

Accordingly, the MPS ordering has become popular in its own right for inequality comparisons in the empirical literature (e.g. Dutta and Foster (2013), Balestra and Ruiz (2015), Madden (2010)). However, these studies draw their conclusions by observing MPS-ordered samples of ordinal variables, without carrying out formal statistical inference. Thus, it is unclear whether the populations underlying these samples are really ordered, or the observed orderings are merely a result of random sampling.

We help to improve on this uncertainty by devising a family of four statistical tests of the null hypothesis that $G$ is *not* an MPS of $F$. We phrase the null hypothesis in this way because researchers are usually interested in the finding that one distribution *is* more unequal than another. By rejecting the hypothesis that $G$ is *not* an MPS of $F$, the researcher is able to conclude that $G$ *is* an MPS of $F$, and therefore that $G$ is more unequal than $F$ in a very robust sense.[2] Each of our four tests is characterised by a test statistic (Likelihood Ratio (LR) or Standardised (Z)), and a method of approximating the distribution of the test statistic under the null (asymptotic or bootstrap). Thus, our family includes both an easy-to-implement test (the asymptotic Z test), and a theoretically-most-attractive test (the bootstrap LR test, see e.g. Davidson and Duclos (2013)), as well as the two intermediate tests.

A natural question arises as to whether and when to prefer the theoretically-most-attractive test over the easy-to-implement test. We answer this by conducting a series of Monte Carlo experiments to compare our tests' performance. We find that all the tests are correctly sized in most cases, although, asymptotic inference can be somewhat oversized when sample sizes are very small or unbalanced. In most other cases, we find little practical difference between the tests in terms of size or power. This finding justifies using the easy-to-implement test in most applications.

Thirdly, we develop new graphical tools for illustrating the null space of the hypothesis of no-ordering, and for comparing the size and power properties of our tests. The null space of our hypothesis is at first hard to visualise because it is defined by the negation of a partial ordering of two ordered multinomial variables. But we show that it can be easily portrayed in the unit square. Moreover, besides being multidimensional and non-convex, the boundary of our null space is the union of two quite different sets: one where the so-called median condition of the MPS fails, and one where the MPS' so-called dominance conditions fail. But by focusing on binomial distributions, we are able to give a representative graphical depiction of the size and power of our tests in and around the boundary.

---

2. See Davidson and Duclos (2013) for a similar framing of the null and alternative hypotheses.

The closest prior research to ours is Gunawan, Griffiths, and Chotikapanich (2018). They conduct Bayesian inference by assuming a uniform prior over all possible distributions and using the observed sample to calculate the respective posterior probabilities that a pair of distributions is or is not ordered. We, on the other hand, conduct frequentist inference which relies exclusively on the data to assess whether there is evidence that the distributions are ordered. Yalonetzky (2013) devises an asymptotic test for first-order stochastic dominance with ordinal variables, which we extend by both testing for equality of the medians, and by reversing the order of the dominance relation above and below the common median. Our work is also related to Abul Naga and Stapenhurst (2015) and Abul Naga, Stapenhurst, and Yalonetzky (2020): while they perform inference on a random variable derived from a particular class of indices consistent with the MPS ordering, we perform inference on the binary outcome given by the partial ordering itself.

The rest of the paper proceeds as follows. Section 3.2 formally defines the MPS partial ordering, motivates our null hypothesis, and introduces a novel graphical representation of the parameter space. Section 3.3 develops our proposed statistical tests. Section 3.4 compares the size and power properties of the four tests with a series of monte carlo experiments. Section 3.5 demonstrate the broad usefulness of our tests in three diverse areas of applications covering subjective wellbeing (happiness in the United States), health economics (self-assessed health in Europe) and development economics (sanitation ladders in Pakistan). Finally section 3.6 offers some concluding remarks.

## 3.2 Median Preserving Spreads and the Null Hypothesis of No Ordering

After presenting the notation, this section defines Allison and Foster (2004)'s MPS partial ordering and describes our null and alternative hypotheses. Then we introduce a novel graphical technique for locating pairs of distributions relative to the null and alternative spaces.

### 3.2.1 Notation

Let $k \in \mathbb{N}$ denote the number of ordered categories and $[k] := \{1, \ldots, k\}$ denote the set of categories. We focus on a pair of samples $(x, y)$ of respective sizes $n_x$ and $n_y$. Each sample is a vector of frequencies which add up to the sample size, for example $x = (x_1, \ldots, x_k) \in \mathbb{N}^k$ and $\sum_{i=1}^{k} x_i = n_x$. Since the states are ordered we can define the cumulants $X = \left( \sum_{j=1}^{1} x_j, \ldots, \sum_{j=1}^{k} x_j \right)$ of $x$; with $Y$ defined analogously for $y$. We use $X_{[i]} := (X_1, \ldots, X_i)$ to denote the first $i$ cumulants of $X$. The sample space is then $\mathcal{X}(k, n_x, n_y) = \{(x, y) \in \mathbb{N}^k \times \mathbb{N}^k \mid X_k = n_x \text{ and } Y_k = n_y\}$. The combined sample is denoted by $W = X + Y$ with combined sample size $n_x + n_y$.

Our ultimate goal is to perform inference on the pair of distributions $(f, g)$ underlying the pair $(x, y)$. Specifically, $x \sim f$ and $y \sim g$ where $f_i$ (respectively $g_i$) denotes the probability that any particular observation from population $f$ (respectively $g$) falls into category $i$. We denote their cumulative distribution functions (henceforth CDF) by $F$ and $G$ respectively. Let $L = (L_1, \ldots, L_k) := W/(n_x + n_y)$ be the sample-weighted average empirical distribution function (EDF).

The parameter space involving all possible pairs of ordered distributions with a given natural number of categories, $k > 1$, is defined by:

$$\Theta := \{(f, g) \in [0, 1]^k \times [0, 1]^k \mid F_k = G_k = 1\}.^3$$

A generic parameter vector is denoted by $\theta = (f, g) \in \Theta$. Samples $x$ and $y$ are drawn from independent ordered multinomial distributions, so the likelihood of $(x, y)$ given $\theta$ is:

$$\mathbb{P}_\theta[x, y] := \frac{n_x!}{\prod_{i=1}^k x_i!} \prod_{i=1}^k f_i^{x_i} \frac{n_y!}{\prod_{i=1}^k y_i!} \prod_{i=1}^k g_i^{y_i}.$$

Also note that $\left(\frac{x}{n_x}, \frac{y}{n_y}\right) \in \Theta$.

Finally, the following sample-size-weighted Euclidean metric for the distance between two distributions will prove useful in several instances below:

$$d(\theta, \theta') = \sqrt{\frac{n_x}{n_x + n_y} \sum_{i=1}^k (f_i - f_i')^2 + \frac{n_y}{n_x + n_y} \sum_{i=1}^k (g_i - g_i')^2}.$$

### 3.2.2   Median Preserving Spreads

We define the median of $F$ to be a category $m \in [k]$ such that $F_{m-1} < 0.5$ and $F_m \geq 0.5$. We assume a unique median category for the purposes of exposition, but all the results generalise to cases with multiple median categories.[4] Allison and Foster (2004) discuss the difficulties of defining suitable measures of dispersion for ordinal variables. They propose the partial ordering over the sample space $\mathcal{X}$ which ranks distributions according to their spread. Here we define the analogous partial ordering for a pair of distributions:

**Definition 2** (Median Preserving Spread). Let $(f, g) \in \Theta$. We say that $g$ is a strict *median preserving spread* (MPS) of $f$, or that $f$ and $g$ are ordered, and write $g > f$, if and only if there exists a category $m$ such that all the following conditions hold:

     [M1] $G_{m-1} < \frac{1}{2}$

---

3.   Even though the sample size $(n_x, n_y)$ is normally considered a parameter of the multinomial distribution, we do not consider it as such because in our applications it is fixed (e.g. by survey design).
4.   For the case of median-preserving spreads with multiple median categories see Kobus (2015).

[M2] $\frac{1}{2} < G_m$

[D1] $G_i > F_i$ for all $i \in [m-1]$

[D2] $G_i < F_i$ for all $i \in [k-1] \setminus [m-1]$ .

We call $f$ the *concentrated* distribution and $g$ the *spread* distribution. If $g$ is not an MPS of $f$ then $f$ and $g$ are unordered, and $g \not\succ f$. If one or more of the strict inequalities holds with equality (which the rest hold strictly) then we say that $g$ is a *weak MPS* of $f$ and write $g \succeq f$. A pair of samples $(x, y) \in \mathcal{X}$ is ordered if and only if the distributions $\frac{x}{n_x}$ and $\frac{y}{n_y}$ are ordered.

### 3.2.3   The Null Hypothesis of No Ordering

We propose tests of the null hypothesis that $g$ is not a strict MPS of $f$ because we are mainly interested in situations where $x$ and $y$ are ordered and want to confirm whether this is indicative of an ordering in the underlying populations. Following Davidson and Duclos (2013), if we reject the null hypothesis that the populations are *not* ordered, then we logically conclude that they *are* ordered.

The *null space* is the subset of all parameter values such that $g$ is not a strict MPS of $f$:

$$\Theta_0 = \{(f, g) \in \Theta \mid g \not\succ f\}.\ ^5$$

Any distributional pair in the null space is denoted by $\theta_0 \in \Theta_0$. The alternative space is the complement of the null space, which is equivalent to the set of all ordered pairs:

$$\Theta_1 := \Theta_0^c = \{(f, g) \in \Theta \mid g \succ f\}$$

A generic element of $\Theta_1$ is denoted by $\theta_1$.

We can graphically depict a two dimensional projection of the null and alternative spaces relative to the whole parameter space. Specifically, a pair of distributions $(f, g) \in \Theta$ can be written as a set of $k$ pairs of cumulants $(F_i, G_i)$. Figure 3.1 plots the pairs of coordinates of distributions $F = (3/24, 9/24, 17/24, 21/24, 1)$ and $G = (7/24, 11/24, 15/24, 17/24, 1)$ in the unit square. Every set of coordinates must include the point $(1, 1)$ and the set of cumulant coordinates is necessarily non-decreasing as we move from left to right because both $F$ and $G$ must be individually non-decreasing. The median category of $F$ (resp. $G$) is given by the state corresponding to the first coordinate to the right of the vertical (resp. horizontal) line at $\frac{1}{2}$. We know the two distributions share the same median category ($m = 3$) because all the coordinates lie in the south-west and north-east quadrants. Any pair of distributions with a coordinate in either the north-west or south-east quadrants do not share the same

---

5.   The set $\Theta_0$ is rotationally symmetric, meaning that reversing the ordering of the categories does not alter the MPS partial ordering of the original distributions. Therefore, all the tests we propose are invariant to reverse ordering of the categories.

median and therefore cannot be ordered. Similarly, we can see that $f$ first-order dominates $g$ below the median and $g$ first-order dominates $f$ at and above the median because all the coordinates lie in the interiors of the two triangles with vertices $(0,0), (\frac{1}{2}, \frac{1}{2}), (0, \frac{1}{2})$ and $(1,1), (\frac{1}{2}, \frac{1}{2}), (1, \frac{1}{2})$, labelled $\Theta_1$ in figure 3.1. It follows from definition 2 that $g$ is an MPS of $f$. In general, $(f, g) \in \Theta_1$ if and only if their coordinates are *all* contained in the triangles labelled $\Theta_1$. Conversely $(f, g) \in \Theta_0$ if and only if *at least one* of their coordinates is contained outside of these triangles, in either of the parallelograms with vertices $(0, \frac{1}{2}), (0, 1), (1, 1), (\frac{1}{2}, \frac{1}{2})$ and $(0, 0), (1, 0), (1, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2})$, labelled $\Theta_0$ in figure 3.1.



**Figure 3.1:** The parameter space, null space and alternative space projected onto the unit square. The red and blue dots illustrate, respectively, the closest distributions on the median and dominance boundaries to the distribution represented by black dots.

The boundary of the null space plays an important role in the proposed tests. Firstly, the Quasi Maximum Likelihood Estimator (QMLE) lies in the boundary of the null space. We will use the QMLE both to calculate the likelihood ratio statistic, and to draw bootstrap samples when carrying out bootstrap inference. Secondly, when we study the empirical size of the tests, we will choose data-generating processes in the boundary of the null space in order to provide an upper bound on the size of the tests of distributions in the interior of the null space. Thirdly, when we study the power of the tests, we will compare the rejection rates of distributions in the alternative space with rejection rates of the corresponding 'closest null' distributions, which lie in the boundary of $\Theta_0$.

We characterise the boundary of the null space as the union of two sets:

**Definition 3.** The *median subset of the boundary* of $\Theta_0$ (henceforth 'median boundary') is the set of all weakly ordered distributions for which at least one of the median constraints in definition 2 hold with equality:

$$\bar{M} = \{(\boldsymbol{F}, \boldsymbol{G}) \in (\Delta[k])^2 \mid \boldsymbol{F} \geq \boldsymbol{G} \text{ and } G_i = \frac{1}{2} \text{ for some } i \in [k]\}.$$

The *dominance subset of the boundary* of $\Theta_0$ (henceforth 'dominance boundary') is the set of all weakly ordered distributions for which at least one of the dominance constraints in definition 2 hold with equality:

$$\bar{D} = \{(\boldsymbol{F}, \boldsymbol{G}) \in (\Delta[k])^2 \mid \boldsymbol{F} \geq \boldsymbol{G} \text{ and } F_i = G_i \text{ for some } i \in [k]\}.$$

The boundary can now be characterised:

**Lemma 6.** *The* boundary *of the null space is equal to the union of the median and dominance boundaries:*

$$\partial\Theta_0 = \bar{M} \cup \bar{D}.$$

*Proof.* See appendix B.1.                                                                                          □

A pair of distributions lies on the median boundary if and only if it has one or more coordinates lying on the horizontal dashed line intersecting the vertical axis at (0,0.5) in figure 3.1, and all other coordinates lying in the interior of $\Theta_0$ triangles. Similarly, a pair of distributions lies on the dominance boundary if and only if it has one or more coordinates lying on the $45°$ dashed line in figure 3.1, and all other coordinates lying in the interior of $\Theta_0$ triangles. We illustrate examples of distributions on the median and dominance boundaries in figure 3.1.

## 3.3   Statistical Tests of No Ordering

A statistical test can be regarded as a function $p : \mathcal{X} \to [0, 1]$ returning a p-value for every sample in the sample space $\mathcal{X}$. The p-value describes the likelihood of observing a sample 'as extreme' as $(\boldsymbol{x}, \boldsymbol{y})$ when the null hypothesis is true, so a low p-value can be taken as evidence that the null hypothesis is false. If the p-value is less than $\alpha \in (0, 1)$ then we 'reject the null hypothesis at the $100\alpha\%$ level.'

A test statistic is a function $\mathcal{S} : \mathcal{X} \to \mathbb{R}$ which formalises what it means for one sample to be 'as extreme' as another by associating each sample with a real number: sample $(\boldsymbol{x}', \boldsymbol{y}')$ is more extreme under the null than $(\boldsymbol{x}, \boldsymbol{y})$ if $\mathcal{S}(\boldsymbol{x}', \boldsymbol{y}') \geq \mathcal{S}(\boldsymbol{x}, \boldsymbol{y})$. Thus we consider four tests of the form $T(\boldsymbol{x}, \boldsymbol{y}) = \mathbb{P}_{\theta_0}[\mathcal{S}(\boldsymbol{x}', \boldsymbol{y}') \geq \mathcal{S}(\boldsymbol{x}, \boldsymbol{y})]$. The remainder of this section discusses our choices of test statistic (LR or Z) and the method of inference (asymptotic or bootstrap).

### 3.3.1  Test Statistics

*The LR Statistic*   The log likelihood ratio (LR) statistic is a natural choice due to its intuitive construction and well-known optimality in terms of uniform power (see section 3.4.2). The LR statistic of a sample $(x, y)$ is the ratio of its unconstrained maximum likelihood function to its constrained counterpart, the QMLE. A large likelihood ratio is evidence that the constraint is hard to satisfy therefore rendering the null hypothesis unlikely to be true.

**Definition 4.** The log likelihood ratio (LR) statistic of a sample $(x, y)$ is given by:

$$\mathrm{LR}(x, y) := 2[\ln(\mathbb{P}_{\theta^*}[x, y]) - \ln(\mathbb{P}_{\tilde{\theta}}[x, y])] \tag{3.1}$$

where $\theta^* \in \arg\max_{\theta \in \Theta} \mathbb{P}_\theta[x, y]$ is the maximum likelihood estimator (MLE), and $\tilde{\theta} \in \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y]$ is the QMLE.

We necessarily have $\mathrm{LR}(x, y) \geq 0$. Lemma 7 derives closed form expressions for the MLE and QMLE.

**Lemma 7.**

1. *The likelihood maximiser (unconstrained) of a sample $(x, y)$ is given by*

$$\theta^* = (x/n_x, y/n_y).$$

2. *If $x$ is not a strict MPS of $y$, then the QMLE is given by*

$$\tilde{\theta} = (x/n_x, y/n_y)$$

   *and $LR(x, y) = 0$.*

3. *Otherwise, if $x$ is a strict MPS of $y$, then the QMLE is given by either one of the following $k - 1$ dominance-constrained distributions $\{\tilde{\theta}^{D_j}\}_{j \in [k-1]} \in \bar{D}$ defined by:*

$$\tilde{\theta}_i^{D_j} = (\tilde{f}_i^{D_j}, \tilde{g}_i^{D_j}) = \begin{cases} \left( \frac{x_i}{X_j} L_j, \frac{y_i}{Y_j} L_j \right) & \text{if } i \leq j \\ \left( \frac{x_i}{n_x - X_j}(1 - L_j), \frac{y_i}{n_y - Y_j}(1 - L_j) \right) & \text{otherwise;} \end{cases}$$

   *or else it is one of the following two median-constrained distributions, $\{\tilde{\theta}^{M_j}\}_{j=m-1,m} \in \bar{M}$ defined by:*

$$\tilde{\theta}_i^{M_j} = (\tilde{f}_i^{M_j}, \tilde{g}_i^{M_j}) = \begin{cases} \left( \frac{x_i}{n_x}, \frac{y_i}{2Y_j} \right) & \text{if } i \leq j \\ \left( \frac{x_i}{n_x}, \frac{y_i}{2(n_y - Y_j)} \right) & \text{otherwise.} \end{cases}$$

   *The likelihood ratio statistic is then given by*

$$LR(x, y) = 2\ln\{ \frac{\mathbb{P}_{\theta^*}[x, y]}{\max\{\mathbb{P}_{\tilde{\theta}^{D_1}}[x, y], ..., \mathbb{P}_{\tilde{\theta}^{D_{k-1}}}[x, y], \mathbb{P}_{\tilde{\theta}^{M_{m-1}}}[x, y], \mathbb{P}_{\tilde{\theta}^{M_m}}[x, y]\}} \}.$$

*Proof.* See appendix B.1 $\qquad\square$

In lemma 7, the multiple cases arise from a Karush-Kuhn-Tucker optimization problem where, depending on the regime of binding constraints, we obtain the various solutions above.

*The Z Statistic* Z statistics have been used in tests of stochastic dominance for multivariate distributions of ordinal variables (e.g. Yalonetzky, 2013). Let $\sigma_i^L = \sqrt{L_i(1-L_i)\frac{n_x+n_y}{n_x n_y}}$ be the standard error of the pooled sample's cumulative frequency $L_i$, and $\sigma_i^Y = \sqrt{(Y_i/n_y)(1-Y_i/n_y)/n_y}$ be the standard error of the sample cumulative frequency $Y_i$. Then consider the Z statistic in definition 5:

**Definition 5.** The Z statistic for a multinomial sample $(\boldsymbol{x}, \boldsymbol{y})$ is given by:

$$Z(\boldsymbol{x}, \boldsymbol{y}) = \min\left\{Z_D^<, Z_D^\geq, Z_M\right\},$$

where $Z_D^< := \min\{(Y_i/n_y - X_i/n_x)/\sigma_i^L \mid i < m_y\}$, $Z_D^\geq := \min\{(X_i/n_x - Y_i/n_y)/\sigma_i^L \mid i \geq m_y\}$ and $Z_M := \min\{(0.5 - Y_{m_x-1}/n_y)/\sigma_{m_x-1}^Y, (Y_{m_x}/n_y - 0.5)/\sigma_{m_x}^Y\}$.

The term $Z_M$ is positive if and only if $m_x = m_y$ (corresponding to conditions [M1] and [M2] in definition 2 for the population counterparts). Hence they are helpful to test the equality of the population medians, which is necessary (but insufficient) to establish an MPS ordering.

The term $Z_D^<$ is the minimum among all the standardised distances of sample cumulative frequencies $\mathbf{Y}$-$\mathbf{X}$ below $m_y$; whereas $Z_D^\geq$ is the minimum among all the standardised distances of sample cumulative frequencies $\mathbf{X}$-$\mathbf{Y}$ at and above $m_y$. Note the similarities with their (unstandardised) population counterparts in conditions [D1] and [D2], respectively. The three statistics are jointly positive, and hence $Z$ is positive, if and only if the sample counterparts of conditions [D1], [D2], [M1] and [M2] hold together. That is, $Z$ is positive *if and only if* $\mathbf{Y} \succ \mathbf{X}$.

### 3.3.2 Method of Inference

The ideal choice of null distribution (Lehmann & Romano, 2005) is that which maximises the probability of the upper contour set $\{\boldsymbol{x}', \boldsymbol{y}' \mid |\mathcal{S}(\boldsymbol{x}', \boldsymbol{y}')| \geq |\mathcal{S}(\boldsymbol{x}, \boldsymbol{y})|\}$, namely $\theta_0 \in \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[|\mathcal{S}(\boldsymbol{x}', \boldsymbol{y}')| \geq |\mathcal{S}(\boldsymbol{x}, \boldsymbol{y})|]$, because this choice ensures that the test always has the correct size (see section 3.4.1). To the best of our knowledge, there is no analytical expression for it in the context of tests involving our specific null space, and numerical solutions are computationally intensive. Instead, we follow the standard approach (e.g. Davidson and Duclos (2013)) of using the QMLE of the observed sample $\theta_0 = \tilde{\theta}$, characterised in lemma 7.

We approximate the probability $\mathbb{P}_{\theta_0}[|\mathcal{S}(\boldsymbol{x}', \boldsymbol{y}')| \geq |\mathcal{S}(\boldsymbol{x}, \boldsymbol{y})|]$ by using either the asymptotic or the bootstrapped distribution of the test statistics. The following theorem will be important for the purpose of asymptotic inference.

**Theorem 3** (Asymptotic distributions of test statistics under the null). *Suppose the true distribution pair lies in the boundary, so that* $(x, y) \sim \theta_0 \in \partial\Theta_0$, *then:*

1. $LR(x, y) \xrightarrow{d} \chi^2(1)$, *and*
2. $Z(x, y) \xrightarrow{d} \mathcal{N}(0, 1)$.

Hence if $\theta_0$ lies in the boundary of the null space, point 1 of the theorem states that the LR statistic converges to a chi-squared variable with one degree of freedom. The one degree of freedom in the chi-squared distribution stems from the difference between the dimensions of the constrained and unconstrained maximum likelihood solutions (Mood, Graybill, & Boes, 1974, p. 440).[6] In the case where $\theta_0$ lies in the interior of the null space, then the LR statistic will generally be lower than for distributions in the boundary, therefore the distribution of the statistic will be first-order stochastically dominated by the $\chi^2(1)$ distribution. We refer the reader to Davidson and Duclos (2013, p. 105). Likewise, the Z statistic is asymptotically standard normal when $\theta_0$ lies in $\partial\Theta_0$, but otherwise is bounded by $\mathcal{N}(0, 1)$. We suggest in practice to approximate the $P$ value of a sample $(x, y)$ by $1 - \zeta[\mathcal{S}(x, y)]$, where $\zeta$ denotes the CDF of $\mathcal{N}(0, 1)$ and $\chi^2(1)$ distributions, respectively. Thus, as we document in our Monte Carlo investigations, the size of the test can be expected to be smaller than the associated nominal value.

Instead of calculating the test statistic of all the samples in the sample space, bootstrap tests approximate the distribution of the test statistic by its empirical distribution in a sample of $B$ samples $\{(x^i, y^i)\}_{i \in [B]}$, each independently drawn from $(x, y)$. The p-value of a sample $(x, y)$ is then approximated by $\#\{(x^i, y^i) \mid \mathcal{S}(x^i, y^i) \geq \mathcal{S}(x, y), i \in [B]\}/B$.

**Theorem 4** (Bootstrap p-values under the null; Davison and Hinkley (1997)). *Suppose the true distribution pair lies in the null space, so that* $(x, y) \sim \theta_0 \in \Theta_0$ *as well as* $(x^i, y^i) \sim \theta_0 \in \Theta_0$ *for all* $i \in [B]$ *where* $B \in \mathbb{N} \setminus \{0\}$, *then the bootstrap p-values for a given statistic* $\mathcal{S}(x, y)$ *are:*

$$T_{BS}(x, y) = \#\left\{(x^i, y^i) \mid |\mathcal{S}(x^i, y^i)| \geq |\mathcal{S}(x, y)|, i \in [B]\right\}/B.$$

Combining the two test statistics with these two methods of approximation gives a family of four tests and respective p-values:

1. Asymptotic Z test: $T_{AZ}(x, y) = 1 - \Phi(Z(x, y))$.
2. Asymptotic LR test: $T_{ALR}(x, y) = 1 - \chi^2(LR(x, y); 1)$.
3. Bootstrap Z test: $T_{BZ}(x, y) = \#\{(x^i, y^i) \mid Z(x^i, y^i) \geq Z(x, y), i \in [B]\}/B$.
4. Bootstrap LR test: $T_{BLR}(x, y) = \#\{(x^i, y^i) \mid LR(x^i, y^i) \geq LR(x, y), i \in [B]\}/B$.

In the next section, we use Monte Carlo simulations to investigate the size and power properties of these four tests.

---

6. See lemma 7 in appendix B.1.

## 3.4  Size and Power Properties

In this section, we introduce novel graphical tools, namely the size-boundary curve for the study of test size, and power-locus curves for the study of test power. We adopt the standard practice of using Monte Carlo experiments to measure the empirical distribution of p-values produced by the tests. Specifically, we draw $M = 100,000$ independent samples $(\boldsymbol{x}^i, \boldsymbol{y}^i)$ from sets of judiciously chosen *data generating processes* (DGPs) $\theta \in \Theta$, and calculate all the p-values, $\{T(\boldsymbol{x}^i, \boldsymbol{y}^i)\}_{i \in [M]}$ for each test $T$. The rejection rate of a nominal size $\alpha$ test at $\theta$ is then estimated by $\#\left\{(\boldsymbol{x}^i, \boldsymbol{y}^i) \mid T(\boldsymbol{x}^i, \boldsymbol{y}^i) \leq \alpha\right\}/M$.

We focus on DGPs with just two categories, i.e. $k = 2$. This class of DGPs is easy to visualise because it is mathematically equivalent to the unit square, with $f_1$ on one axis and $g_1$ on the other. The median boundary is mathematically equivalent to the horizontal line intersecting the vertical axis at $(0, 0.5)$ and the dominance boundary is mathematically equivalent to the 45° line. Because the boundary is unidimensional it is easy to show how rejection rates vary along it. Similarly, we are able to identify a unidimensional 'interior locus' which allows us to illustrate how power varies against different DGPs in the alternative space. We argue in section 3.4.3 that tests' behaviour in the $k = 2$ case is indicative of behaviour in higher dimensions.

### 3.4.1  Size

If the *empirical* size, $\mathbb{P}_{\theta_0}[T(\boldsymbol{x}, \boldsymbol{y}) < \alpha]$, of a test $T$ is less than the *nominal* size $\alpha$ for all null distributions $\theta \in \Theta_0$, then we say that the test is *correctly sized* at level $\alpha$; otherwise it is *oversized*. A standard size curve plots the actual (empirical) rejection rate of a test against its nominal size for a given null distribution. However, the intricacies of our null space are difficult to capture with a small selection of null distributions. Instead, we build a comprehensive picture of behaviour in the boundary, by holding the nominal size constant at the 5% level and plotting the empirical size of our tests against a grid of different DGPs in the boundary of the $k = 2$ null space, for a range of sample sizes. Our tests will have higher rejection rates on the boundary than anywhere else in the null space, so this procedure gives an upper bound on the actual size of the tests. Figure 3.2 illustrates the DGP's used for the cases $n_x = n_y = 10, 100, 1000.$[7] Our interest in small sample sizes, and more specifically in small ratios of $n_x$ to $n_y$ and $n_y$ to $n_x$ is three-fold: (a) to investigate the relative merits of the bootstrap versus asymptotic inference in relation to the size and power of Z and LR tests; (b) to explore lower bounds on sample size in relation to the performance of the tests; and (c) to highlight the asymmetric role of sample sizes of the spread distribution ($n_y$) and the concentrated distribution ($n_x$) in the statistical performance of the four tests.

---

7.  The precise choices of boundary DGPs for these and other sample sizes are listed in appendix B.2.

**Figure 3.2:** Boundary DGPs used to generate size curves for cases $k = 2$ and $n_x = n_y$.

*Results* Figure 3.3 shows the rejection rate of all four tests at the 5% nominal level, as a function of the first coordinate of the boundary DGPs. In each panel, the first half of the horizontal axis, from 0 to 0.5, corresponds to the median boundary (moving along the horizontal dotted line from coordinate (0,0.5) to (0.5,0.5) in figure 3.2); and the second half of the horizontal axis, from 0.5 to 1, corresponds to the dominance boundary (moving along the diagonal dotted line from coordinate (0.5,0.5) to the origin in figure 3.2). The intersection of the median and dominance boundaries coincides with the point 0.5 (the kink of the dotted line in the middle of figure 3.2). From top-left downward and rightward, panels in row $i$ show results for $n_x = 10^i$ while panels in column $j$ show results for $n_y = 10^j$, where $j = 1, 2, 3$.

All the tests are correctly sized in most cases. Exceptions arise on the dominance boundary when $n_x$, the size of the sample drawn from the more concentrated distribution, is small relative to the size of the sample drawn from the more polarised distribution. In these cases, tests based on the Z statistic can have sizes more than double their nominal levels. Meanwhile, the sizes of tests based on the LR statistic are never more than 20% above their nominal values. Two other exceptions occur: the asymptotic tests are both slightly oversized near the end of the median boundary in the case $(n_x, n_y) = (10, 1000)$, and the bootstrap LR test is slightly oversized near the end of the dominance boundary in the case $(n_x, n_y) = (1000, 10)$.

The rejection rates of all tests drop to zero near the intersection of the median and dominance boundaries, especially when the sample size of the less concentrated distribution is small. The region of the boundary where the rejection rate drops to zero vanishes as the sample sizes increase. The lower rejection rates vis-a-vis those in other points in the boundary are not surprising: for points other than (0.5,0.5), the proportion of neighbouring distributions that belong to the null and alternative space are of equal size, namely 1/2 and 1/2. However,

**Figure 3.3:** Size-boundary curves (for nominal 5% tests). Key: Solid/light blue — bootstrap LR; dashed/light red — bootstrap Z; dotdash/dark blue — asymptotic LR; dotted/dark red — asymptotic Z.

at (0.5,0.5), the proportion of neighbouring distributions that belong to the null space is now equal to 3/4, whereas the proportion of neighbouring distributions that belong to the alternative space is now equal to 1/4. For this reason, the probability of a sample with an empirical distribution in the null space is more likely, leading to fewer rejections of the null hypothesis.

### 3.4.2 Power

Besides correct size, the other crucial property for statistical tests is the ability to distinguish between true and false hypotheses, known as the 'power' of the test. Davidson and MacKinnon (1998) propose to assess power by plotting 'size-adjusted' size-power curves. These curves are constructed by plotting the rejection rate for a distribution in the alternative space against the rejection rate for the closest corresponding distribution in the null space. In the case $k = 2$, there are only two candidates for the closest null distribution: the closest

distribution in the median boundary, $\theta^M := \left(f, \left(\frac{1}{2}, \frac{1}{2}\right)\right)$; and the closest distribution in the dominance boundary, $\theta^D := \left(\frac{n_x f + n_y g}{n_x + n_y}, \frac{n_x f + n_y g}{n_x + n_y}\right)$.[8] Figure 3.1 illustrates both the closest pair of distributions on the median boundary, denoted by red circles, and the closest pair of distributions on the dominance boundary, illustrated by blue circles, to the pair of distributions denoted by black circles, for $k = 2$.

As with size, we face a choice over which alternative distributions to study. Because the boundary separating the null and alternative spaces of the tests introduced in this paper arises as the union of the dominance and median boundaries, we focus on studying power against alternatives that are equidistant from these median and dominance boundaries. We refer to these DGPs as the 'interior locus'. Figure 3.4 shows the grid of DGPs on the interior locus, connected by a solid blue line, that we use for our experiments with $n_x = n_y$.[9] We also show, for each of these alternative DGPs, the two closest null distributions — one on each boundary — connected to the interior locus by a red dashed line. This interior locus is worth studying because it partitions the alternative space into a set of DGPs which are closest to the median boundary and a set of DGPs closest to the dominance boundary, and every DGP in the alternative space can be uniquely identified with a point on the interior locus which shares the same closest null distribution (be it on the median or the dominance boundary). Moreover, we expect all the tests to have lower power against an arbitrary alternative DGP than against its counterpart in the interior locus. Thus, the interior locus provides an upper bound on the test's power.



**Figure 3.4:** Alternative DGPs used to generate power curves for cases $k = 2$ and $n_x = n_y$.

---

8. We measure distance with the Euclidean metric defined in subsection 3.2.1.
9. Precise values for these DGPs, and those used for other ratios of $n_x$ to $n_y$ are listed in appendix B.2.

*Results*    In figure 3.5 we introduce a novel *power-locus curve* used to investigate the power properties of the various tests. By definition, the alternative DGPs in the 'interior locus' have two closest null distributions. Therefore there are two ways to evaluate power against these DGPs. Each panel of figure 3.5 illustrates both methods for a different pair of sample sizes. The first half of the horizontal axis depicts the 'median power curve': the power against each alternative DGP from left to right in terms of figure 3.4, calculated using the closest null on the *median* boundary. The second half depicts the reflected 'dominance power curve': power against each alternative calculated using the closest null on the *dominance* boundary and in the reverse order. Such a display of the results enables us to see how the power varies as the DGP approaches the intersection of the two boundary lines from the 'median direction', and from the 'dominance' direction, respectively. We expect that power against alternatives near the median boundary will behave similarly to the median power curve, and that power against alternatives near the dominance boundary will behave similarly to the dominance power curve.
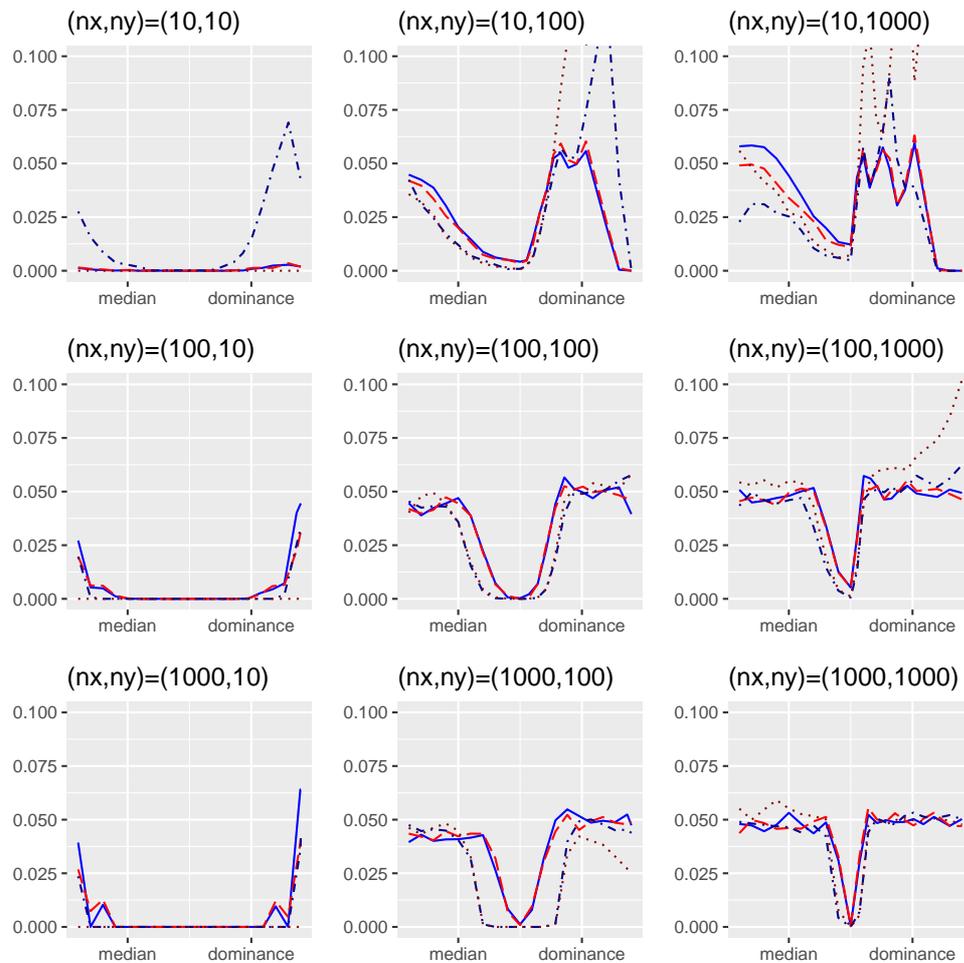


**Figure 3.5:** Power-locus curves (for nominal 5% tests). Key: Solid/light blue — bootstrap LR; dashed/light red — bootstrap Z; dotdash/dark blue — asymptotic LR; dotted/dark red — asymptotic Z.

All the tests can perfectly distinguish some alternatives from the null space whenever both sample sizes are above 100. By the time both sample sizes reach 1000, the tests can perfectly distinguish a false hypothesis along distributions pertaining to three quarters of the interior locus. However, power drops rapidly when the sample size falls below 100: power halves when $n_x$ falls from 100 to 10, and reduces by a factor of 4 when $n_y$ falls from 100 to 10. All the tests have more or less the same power when sample sizes are both in the order 100 or higher. Even with smaller sample sizes, the choice of test statistic appears to have minimal impact on power. However, there is evidence of systematic disparities between asymptotic and bootstrap inference for smaller sample sizes. When $n_x$ is both small relative to $n_y$ and very small in absolute terms, the asymptotic tests are more powerful against some distributions nearer the median boundary. However, when $n_y$ is very small, the power of asymptotic inference is both erratic and consistently lower than the power of bootstrap inference.

### 3.4.3 More than Two Categories

A distribution $(f, g)$ with $k > 2$ categories can be decomposed into $k - 1$ two-category distributions $(f^i, g^i)$ defined by $(F^i, G^i) = ((F_i, 1), (G_i, 1))$ for any $i \in [k - 1]$. The rejection rate of a test at $(f, g)$ is therefore a function of the $k - 1$ rejection rates at each of the $(f^i, g^i)$. For example, intersection-union tests (Berger, 1982) reject the null that $(f, g)$ are not ordered if and only if they reject all the $k - 1$ hypotheses that each of the $(f^i, g^i)$ are not ordered. Graphically, this means that we infer that all the coordinates in figure 3.1 are contained within the two triangles representing the alternative space, if and only if we infer that the coordinate closest to the edge of the triangles is inside. The study of real world DGP's characterised by more than two categories is taken up in section 3.5.

## 3.5 Empirical Illustrations

We consider three real-world inequality assessments: happiness in the United States, self-reported health in a set of European countries, and sanitation ladders in Pakistan. We undertake 499 bootstrap replications in each of the three applications. In the context of self-reported health and sanitation ladders, we present p-value curves constructed from 10,000 Monte Carlo simulations.

### 3.5.1 Happiness Inequality in the United States

We revisit the study of Dutta and Foster (2013) on happiness inequality in the United States. They use data from the U.S. General Social Survey (GSS) between 1972 and 2010 (Dutta & Foster, 2013, table 1, p. 402). The GSS asks the following ordered-response question on wellbeing: "Taken all together, how would you say things are these days — would you say that you are 'very happy', 'pretty happy' or 'not too happy?'" Dutta and Foster (2013) did not test whether the documented ordering of happiness distributions was statistically significant. The family of tests developed in this paper provides the required statistical inference.

Table 3.1 reports p-values of the bootstrap LR test, where an entry in row $i$ of column $j$ is the p-value of the sample under the null hypothesis that year $j$ distribution is not an MPS of year $i$ distribution. A blank cell indicates that column $j$ sample is not an MPS of the row $i$ sample. Our results show that most of these inequality comparisons are (individually) statistically significant, with p-values close to 0. The inferential exercise broadly supports the underlying pattern identified by Dutta and Foster (2013), namely a fall in happiness inequality across the 70s and 80s, that is reversed in the 90s and 2000s.

There are, however, some noteworthy exceptions. For example, consider the finding that 'happiness inequality was lower in 1985 compared with seventeen other years' (Dutta & Foster, 2013, p. 405). Table 3.1 reveals that the p-value of the comparison between 1985 and 1998 equals 0.11. Other comparisons where the p-value is 10% or higher include the (1993, 2004) and (1993, 2006) pairs (both with a p-value of 0.14), the (1998, 1975) pair (p-value of 0.16), the (2000, 1987) pair (p-value of 0.27), the (2004, 1982) pair (p-value of 0.50), and the (2006, 2004) pair (p-value of 0.40). If one adopts the standard convention of failing to reject a null hypothesis when the associated p-value exceeds the 5% level, one would not find sufficient statistical evidence to support the conclusion that distributions $i$ and $j$ were ordered in the aforementioned comparisons.

### 3.5.2 Inequality in Self-assessed Health in the European Union

Self-assessed health (SAH) measures are increasingly used in health surveys, as such subjective assessments of well-being have shown to be strong predictors of morbidity as well as mortality (Latham & Peek, 2013). The Survey on Incomes and Living Conditions (SILC) conducted by EUROSTAT collects data on five levels of self-assessed health in the European Union. Respondents choose from the following ordered subjective health categories: (1) very bad, (2) bad, (3) fair, (4) good, and (5) very good. In 2017, the multinomial distributions of the Netherlands and Denmark were, respectively, $x/n_x = (0.01, 0.04, 0.19, 0.54, 0.22)$ with a sample size $n_x = 13328$ and $y/n_y = (0.03, 0.06, 0.21, 0.45, 0.25)$ with a sample size $n_y = 5906$. The two samples are ordered: sharing 'good health' as the median category and

**Table 3.1:** Bootstrap LR $p$-values for Dutta and Foster (2013, table 2).

| | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 80 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 93 | 94 | 96 | 98 | 00 | 02 | 04 | 06 | 08 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 08 | | | | | | | | | .02 | | | | | | | | | | | | | | .00 | | | | |
| 06 | | | | | | | | | .00 | | | | | .00 | .14 | | | | | | | | | | | .27 | |
| 04 | | | | | | | | | .00 | | | | | .00 | .14 | .00 | .00 | .02 | | | .00 | | | | | .40 | |
| 02 | | | | | | | | | | | | | | | .02 | | .00 | | | | | | | | | | |
| 00 | | | | | | | | | .00 | | | | | | | .05 | | | | | | | | | | | |
| 98 | | | | | | | | | .11 | | | | | | | .00 | | .00 | | | | | | | | | |
| 96 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 94 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 93 | | | | | | | | | .02 | | | | | | | .01 | | | | | | | | | | | |
| 91 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 90 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 89 | | | | | | | | | .02 | | | | | | | | | | | | | | | | | | |
| 88 | | | | | | | | | | | | | | | .03 | | | | | | | | | | | | |
| 87 | | | | | | .00 | | | .00 | | | | .00 | .00 | .00 | .00 | .02 | | .27 | | | | | | | | |
| 86 | | | | | | | | | | | | | .00 | | .00 | | | | | | | | | | | | |
| 85 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 84 | | | | | .00 | .00 | | | .00 | | .09 | .00 | .00 | .00 | .00 | .00 | .00 | | | | .00 | | | | | | |
| 83 | | | | | | | | | | | | .00 | | | | | | | | | | | | | | | |
| 82 | | | | | | | | | .00 | | .03 | | .00 | .00 | .00 | .00 | .00 | .50 | .00 | | | | | | | | |
| 80 | | | | | .00 | .00 | | | .00 | .03 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | | .00 | | | | | | | |
| 78 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 77 | | | | .04 | | | | | .00 | .00 | | .00 | .45 | .00 | | .00 | .10 | .00 | | .00 | | | | | | | |
| 76 | | | .44 | .09 | | | | | .00 | | .02 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | | | | | | | | | |
| 75 | | | | | | | .00 | .02 | | | .00 | | | .00 | .16 | | .00 | | | | | | | | | | |
| 74 | .03 | | .01 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | | |
| 73 | | | .38 | .00 | .00 | .01 | .00 | .05 | .00 | .00 | | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | | | | | |
| 72 | .00 | | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .05 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | .00 | | | |

with Denmark's distribution being a MPS of the Netherlands'. As is typical of distributions of self-assessed health, both distributions exhibit some class imbalances, with near-zero probability mass associated with the bottom health categories, and with over 40% mass attached to the median category.

We investigate a Z-test and a likelihood ratio test of the null hypothesis that the Danish distribution is not a MPS of the Dutch distribution. Figure 3.6 shows that all the tests can perfectly distinguish this false hypothesis from the closest (true) null hypothesis (the power curves are vertical). Next, we turn our attention to comparing the four tests in terms of the p-value curves in the figure. For all relevant nominal test sizes (0 to 10%), all four tests are correctly sized: the asymptotic Z test is undersized, while the actual size of the other

three tests coincides with the nominal size. In the context of this application, pertaining to large samples associated with distributions of self-assessed health, it is not possible to infer whether the Z or LR test is preferable in terms of size. Furthermore, this conclusion remains unchanged when we either consider the asymptotic or bootstrap approximation of the related test statistics.
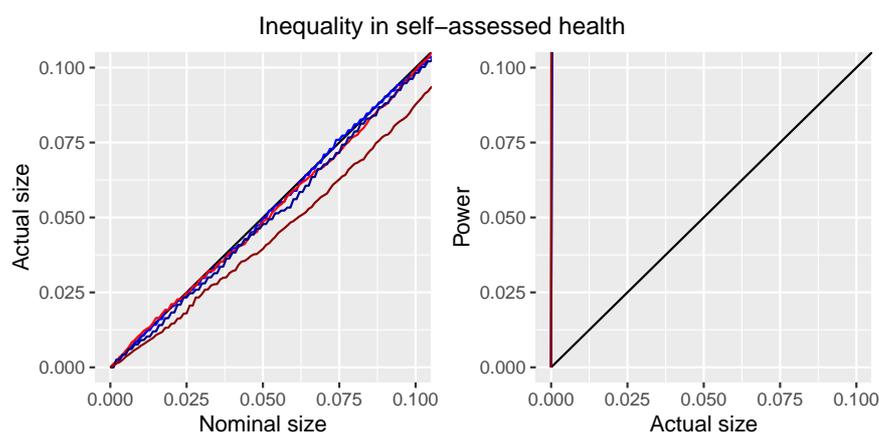


**Figure 3.6:** Left: Size curves for the distribution $(F, G) = ((0.02, 0.06, 0.24, 0.78, 1.00), (0.02, 0.08, 0.29, 0.75, 1.00))$. Right: Power curves for the distribution $(F, G)=((0.01, 0.05, 0.24, 0.78, 1.00), (0.03, 0.09, 0.30, 0.75, 1.00))$.

### 3.5.3 Inequality in Sanitation Ladders in Pakistan

Improvements in access to toilet facilities are a measure of living standards related to sanitation in developing countries (Seth & Yalonetzky, 2021). The 2017-8 Demographic and Health Survey of Pakistan collects data on different forms of sanitation facilities which can be grouped into a four-level sanitation ladder following the guidelines of the Joint Monitoring Program by the WHO and UNICEF.[10] The ensuing ordered categories are the following: (1) open defecation, (2) access to an unimproved toilet facility (buckets and latrine toilets that do not flush), (3) shared improved toilet facilities (such as a shared toilet that flushes to piped sewer system) and finally (4) improved toilet facility that is not shared. The probability mass distributions pertaining to Islamabad (the capital city) and Baluchistan are, respectively, $x/n_x = (0.003, 0.001, 0.060, 0.936)$ with a sample size $n_x = 1295$ and $y/n_y = (0.135, 0.039, 0.142, 0.684)$ with a sample size $n_y = 1521$. The two distributions highlight important regional differences in attainment in sanitation, and furthermore in spread. That is, Baluchistan's sample is a MPS of Islamabad's.

---

10. See https://washdata.org/.

There are three interesting properties these distributions exhibit in terms of statistical inference. Firstly, the median state in both distributions is $m = 4 = k$ (the top category). Were it not the case that our proposed tests jointly test that the distributions share an equal median, and are ordered according to the MPS criterion, the inferential exercise here would be equivalent to a test of first-order stochastic dominance (of Islamabad over Baluchistan). However, as we have earlier emphasised, the critical region of the tests in this paper is constrained by the union of the dominance and median boundaries, and in this sense, a simple test of first-order dominance would not provide a valid inferential tool in this application. The second interesting property of the data is that both distributions exhibit severe class imbalances, with the highest sanitation state (improved toilet facility that is not shared) being associated with a probability mass in excess of 66%. Finally, the Islamabad distribution has a probability mass of 0.001 (one unique observation) in the second sanitary ladder state.

We investigate a Z-test and an LR test of the null hypothesis that the Baluchistan distribution is not a MPS of Islamabad. As in the EU health application, all the tests can perfectly distinguish this false hypothesis from the closest (true) null hypothesis, therefore they reproduce the power curves.[11] We, therefore, focus our attention on comparing the four tests in terms of the p-value curves, plots of which are provided in figure 3.7. The dark red curve presents the p-value plot for the asymptotic approximation of the Z test, while the light red curve pertains to the asymptotic LR test. The dark blue curve is the p-value curve corresponding to the bootstrap approximation of Z statistic's sampling distribution, while finally, the light blue curve refers to the bootstrap LR test.
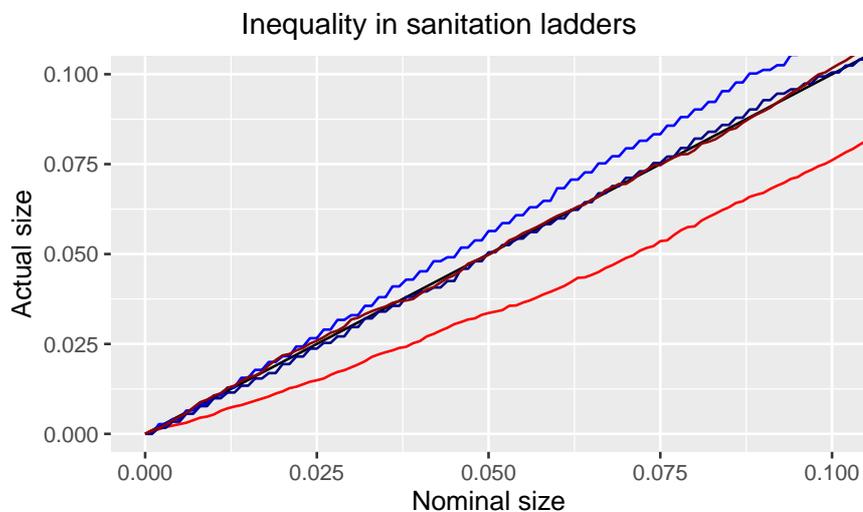


**Figure 3.7:** Size curves for the distribution $(\boldsymbol{F}, \boldsymbol{G})$=((0.07, 0.07, 0.13, 1.00), (0.07, 0.12, 0.27, 1.00)).

---

11. They are available upon request.

We summarise our findings for all relevant nominal test sizes (0 to 10%) as follows: the asymptotic LR test is correctly sized (the relevant p-value curve is below the 45-degree line), the asymptotic Z and the bootstrap Z tests have an actual size equal to their nominal size (the p-value curves are on the 45-degree line). On the other hand, the bootstrap LR test is moderately oversized (the p-value curve lies above the 45-degree line). Overall though, all four tests perform satisfactorily in the context of this application, with asymptotic inference being preferable over bootstrap inference. We attribute the better performance of asymptotic inference in the context of this application to the occurrence of severe class imbalances in both distributions.

## 3.6 Conclusion

The purpose of this paper was to introduce a family of tests for the hypothesis that an ordered multinomial distribution $G$ is a MPS of $F$. Using Monte Carlo simulations, we found that the choice between Z and LR test statistics does not have a large impact on the tests' properties, but the method used to approximate the sampling distribution of the statistics under the null does. In a wide range of data generating processes, bootstrap inference generally exhibited better size and power properties than asymptotic inference. We have further illustrated the proposed tests in three areas of inequality applications: happiness in the United States, self-assessed health in Europe and sanitation ladders in Pakistan.

The paper can be extended in several directions. MPS is a special case of Mendelson (1987)'s "quantile-preserving spread" partial ordering: a distribution $G$ is a *quantile-preserving spread* (QPS) of a distribution $F$ ('$F$ and $G$ are ordered') if there exists a category $q$ such that the mass of $G$ lies further away from the state $q$ than does the mass of $F$ ($G$ has a thicker tail than $F$ around $q$). If this is true for $q$ equal to the lowest (respectively highest) category then $G$ first order dominates $F$ (respectively $F$ first order dominates $G$), so QPS orders the distributions according to their relative locations. In intermediate cases where $q$ is equal to neither the highest nor the lowest category then QPS is sensitive to both the location and variability of the distributions. Tests of quantile preserving spreads can be formulated by replacing the median boundary with an appropriate quantile boundary.

Likewise, one may derive tests of hypotheses constructed from linear transformations of the vector of contrasts related to the median preserving spreads ordering; for instance, the bipolarization partial order of S. Chakravarty and Maharaj (2012). Finally, we mention the need to develop exact inference for tests of median-preserving spreads, yielding the p-values of every conceivable sample, as a companion method to the bootstrap and asymptotic methods of inference introduced in this paper.

# Lemons by Design

## A.1 Mathematical appendix

### A.1.1 Proof of proposition 8

The cost of the scheme $\mathcal{S}^{II} = (q^{II}, (f_e^{II}, f_0^{II}, r^{II}))$ is

$$\mathbb{E}[\pi_c(r^{II}(x_I) - f_e^{II}(x_I)) + (1 - \pi_c)(-f_0^{II}(x_I))]$$

$$= \int_{\frac{k-\Pi}{k}}^{1} \pi_c r^{II}(x_I) - (\pi_c f_e^{II}(x_I) + (1 - \pi_c)f_0^{II}(x_I))dx_I$$

$$= \int_{\frac{k-\Pi}{k}}^{1} \pi_c\left(k - \frac{k-\Pi}{x_I}\right) - (\pi_c(1 - \pi_c)k - (1 - \pi_c)\pi_c k)dx_I$$

$$= \pi_c \int_{\frac{k-\Pi}{k}}^{1} k - \frac{k-\Pi}{x_I}dx_I$$

$$= \pi_c\left(k\left(1 - \frac{k-\Pi}{k}\right) - (k - \Pi)\left[\ln x_I\right]_{\frac{k-\Pi}{k}}^{1}\right)$$

$$= \pi_c\left[\Pi + (k - \Pi)\ln\left(1 - \frac{\Pi}{k}\right)\right]$$

$$\leq \pi_c\left[\Pi - (k - \Pi)\frac{\Pi}{k}\right]$$

$$= \pi_c \frac{\Pi^2}{k} \xrightarrow{k\to\infty} 0,$$

where the inequality comes from the fact that $\ln(1 + x) \leq x$ for all $x > -1$ (Topsøe, 2004).

We now show that the scheme is optimal by showing that any feasible solution to the government's problem costs weakly more than $\mathcal{S}^{II}$. In any scheme where Ina has full information about the transfers, she will accept the bribe $b$ whenever she receives a message $x_I$ such that $r(x_I) < b$. Finn anticipates this, so he accepts bribe $b$ if $\mathbb{E}[f_e(x_I) - f_0(x_I)|r(x_I) < b] > b$. Therefore, any informed inspector scheme which deters bribes must at least satisfy $\mathbb{E}[f_e(x_I) - f_0(x_I)|r(x_I) < b] \leq b$ for all bribes $b \geq 0$.

Define Ina's expected incentive when Ina rejects a bribe $b$ by $E(b) := \mathop{\mathbb{E}}\limits_{x \sim q}\left[f_e(x_I) - f_0(x_I) | r(x_I) \geq b\right]$, and $F_r(b) = \mathop{\mathbb{P}}\limits_{x \sim q}\left[r(x_I) < b\right]$. We use the fact that $E(0) = \mathbb{E}[f_e(x_I) - f_0(x_I)] = \mathbb{E}[f_e(x_I) - f_0(x_I) | r(x_I) < b] F_r(b) + E(b)(1 - F_r(b))$ for any $b$, to rewrite the no-bribery constraint as $\frac{E(0) - E(b)(1 - F_r(b))}{F_r(b)} \leq b$. This rearranges to give $F_r(b) \leq \frac{E(b) - E(0)}{E(b) - b}$. If $E(0) < b$ then $\frac{E(b) - E(0)}{E(b) - b} > 1$, so the constraint is slack. Otherwise, if $E(0) \geq b$, then $\frac{E(b) - E(0)}{E(b) - b}$ is increasing in $E(b)$ and decreasing in $E(0)$. Finn's limited liability constraint requires that $E(b) \leq k$ and his incentive compatibility constrain requires that $E(0) \geq \Pi$. It follows that every bribery-proof informed inspector scheme satisfies

$$F_r(b) \leq \frac{E(b) - E(0)}{E(b) - b} \leq \frac{k - \Pi}{k - b}.$$

In $\mathcal{S}^{II}$, the distribution of rewards is

$$F_r^{II}(b) = \mathop{\mathbb{P}}\limits_{x_I \sim q^{II}}\left[r^{II}(x) < b\right] = \mathop{\mathbb{P}}\limits_{x_I \sim q^{II}}\left[x_I < \frac{k - \Pi}{k - b}\right] = \frac{k - \Pi}{k - b}.$$

Hence the distribution of rewards in any bribery-proof solution must first order stochastically dominate the distribution of rewards in the scheme $\mathcal{S}^{II}$. This implies that every bribery-proof solution has a weakly higher expected reward than does $\mathcal{S}^{II}$. At the same time, $\mathcal{S}^{II}$ exactly satisfies Finn's (VP) constraint, so every feasible scheme must have a weakly lower expected fine than $\mathcal{S}^{II}$. The cost of a scheme is given by Ina's expected reward minus Finn's expected fine, so it follows that every feasible scheme must cost weakly more than $\mathcal{S}^{II}$.

### A.1.2   Informed firm schemes

Consider the biased coin toss (informed firm) scheme described in table A.1. Similar tech-

**Table A.1:** The biased coin toss (informed firm) scheme.

| | Lemon | Peach |
|---|---|---|
| Probability | $1 - \sqrt{1 - \Pi/\kappa}$ | $\sqrt{1 - \Pi/\kappa}$ |
| Fine Finn | 4 | $\left(1 - \sqrt{1 - \Pi/\kappa}\right)\kappa$ |
| Reward Ina | 4 | 0 |

niques to those used in section 2.4.1 show that this scheme satisfies voluntary participation, incentive compatibility and deters bribes. The cost of the scheme is $\pi_c \sqrt{\kappa}(\sqrt{\kappa} - \sqrt{\kappa - \Pi})$. When $\kappa = 4$ and $\Pi = 3.9$ we get that this informed firm scheme costs $3.4\pi_c$, whereas the cheapest informed inspector scheme costs $\pi_c\left[\Pi + (k - \Pi)\ln\left(1 - \frac{\Pi}{k}\right)\right] = 3.5\pi_c$. In general, informed firm schemes are cheaper when $\kappa$ is small enough relative to $\Pi$.

The optimal informed firm scheme takes the following form,

$$q^{\mathsf{IF}} \text{ is uniform on } [0,1]$$

$$f_e^{\mathsf{IF}}(x_F) = (1 - \pi_c) \begin{cases} k & \text{if } x_F \le \tilde{x}^{IF} \\ \dfrac{\tilde{x}^{IF}}{x_F} & \text{otherwise} \end{cases}$$

$$f_0^{\mathsf{IF}}(x_F) = -\pi_c \begin{cases} k & \text{if } x_F \le \tilde{x}^{IF} \\ \dfrac{\tilde{x}^{IF}}{x_F} & \text{otherwise} \end{cases}$$

$$r^{\mathsf{IF}}(x_F) = \begin{cases} k & \text{if } x_F \le \tilde{x}^{IF} \\ 0 & \text{otherwise,} \end{cases}$$

where $\tilde{x}^{IF}$ solves $\tilde{x}^{IF} \ln\left(\frac{e}{\tilde{x}^{IF}}\right) = \frac{\Pi}{\kappa}$. Similar techniques to those used in appendix A.1.1 show that this scheme is feasible and cheaper than any other informed firm scheme. However, it is not easy to work with because no analytical solution for $\tilde{x}^{IF}$ exists.

### A.1.3   Collusion proofness

In the main body of the paper, we restrict attention to bribery contracts. The revelation principle tells us that for every equilibrium of every bribery game, there exists an incentive compatible collusive side contract. An incentive compatible collusive side contract is a generalisation of a bribery contract in which the players can play correlated strategies and in which side of the bribe can depend on the players' types. We use this more general notion of collusion to prove that the government can restrict attention to collusion-proof contracts, and therefore to bribery-proof contracts.

Our notion of collusive side contracts is inspired by Laffont and Martimort (1997). Suppose that a fourth player, Marta the Mafia, offers to enforce collusive side contracts. In a direct collusive side contract, Ina and Finn report their private messages $x = (x_I, x_F)$ to Marta; Marta tells Ina to be silent with probability $s(x)$; and Marta makes (potentially negative) transfers $b_F(x)$ to Finn and $b_I(x)$ to Ina. We denote a collusive side contract by $\mathcal{C} = (s, b_F, b_I)$. A side contract $\mathcal{C}$ is budget balanced if Marta makes does not lose money on average, i.e. $\mathbb{E}_x[b_F(x) + b_I(x)] \le 0$. It is incentive compatible if it satisfies the usual incentive compatibility constraints that ensure it is in Ina and Finn's best interest to report their message truthfully. Finally, Marta cannot force Ina and Finn to participate, so $\mathcal{C}$ must satisfy the usual voluntary participation constraints. A balanced budget, incentive compatible, voluntary side contract is *interim efficient* if it delivers a strictly higher payoff to at least one type of one player, and a weakly higher payoff to all types of both players, than every other balanced budget, incentive compatible, voluntary side

contract. A scheme $(q, (f_e, f_0, w))$ is *weakly collusion-proof* if no interim efficient side contract exists. Hence, weak collusion-proofness implies bribery-proofness, but the reverse does not necessarily hold. In the main paper, we show that our scheme is bribery-proof. We conjecture that it is also weakly collusion-proof.

### A.1.4 The two sided scheme costs strictly less than the one sided schemes.

The informed inspector scheme costs strictly more than the two-sided scheme at all parameter values:

$$
\begin{aligned}
c^{\text{II}}/\pi_c &= \Pi + (\kappa - \Pi)\ln\left(1 - \frac{\Pi}{\kappa}\right) \\
&\geq \Pi - (\kappa - \Pi)\frac{\frac{\Pi}{\kappa}}{\sqrt{1 - \frac{\Pi}{\kappa}}} \\
&= \Pi - (\kappa - \Pi)\frac{\frac{\Pi}{\sqrt{\kappa}}}{\sqrt{\kappa - \Pi}} \\
&= \Pi - \sqrt{\kappa - \Pi}\frac{\Pi}{\sqrt{\kappa}} \\
&= \frac{\Pi}{\sqrt{\kappa}}\left(\sqrt{\kappa} - \sqrt{\kappa - \Pi}\right) \\
&> (\sqrt{k} - \sqrt{\kappa - \Pi})^2 \\
&= c^*/\pi_c,
\end{aligned}
$$

where the first inequality results from the fact that $\ln(1 + x) \geq \frac{x}{\sqrt{1+x}}$ for all $x \in (-1, 1]$ (Topsøe, 2004), and the second comes from the fact that

$$
\begin{aligned}
\sqrt{\kappa} &> \sqrt{\kappa - \Pi} \\
\sqrt{\kappa}\sqrt{\kappa - \Pi} &> \kappa - \Pi \\
\Pi &> \kappa - \sqrt{\kappa}\sqrt{\kappa - \Pi} \\
\frac{\Pi}{\sqrt{\kappa}} &> \sqrt{\kappa} - \sqrt{\kappa - \Pi}.
\end{aligned}
$$

Now we show that the two-sided scheme costs half as much as the informed inspector scheme in the limit.

$$
\begin{aligned}
\frac{c^*}{c^{ll}} &\leq \frac{\sqrt{\kappa}}{\Pi}(\sqrt{\kappa} - \sqrt{\kappa - \Pi}) \\
&= \frac{\kappa - \sqrt{\kappa(\kappa - \Pi)}}{\Pi} \\
&= \frac{\frac{\kappa - \sqrt{\kappa(\kappa - \Pi)}}{\Pi}(\kappa + \sqrt{\kappa(\kappa - \Pi)})}{\kappa + \sqrt{\kappa(\kappa - \Pi)}} \\
&= \frac{\frac{\kappa^2 - \kappa(\kappa - \Pi)}{\Pi}}{\kappa + \sqrt{\kappa(\kappa - \Pi)}} \\
&= \frac{\kappa(\kappa - \kappa(\kappa - \Pi))}{\Pi} \\
&= \frac{\kappa}{\kappa + \sqrt{\kappa(\kappa - \Pi)}} \xrightarrow{\kappa \to \infty} \frac{1}{2},
\end{aligned}
$$

where the first inequality comes from the previous calculations.

### A.1.5 Derivation of Conditional Expectations

Here we derive equation (2.2). The derivation of equation (2.3) is completely analogous.

$$
\begin{aligned}
\mathbb{E}\left[\min\{1, x_F^\lambda x_I^{-1}\} \mid x_I, x_F \leq y_F\right] &= \begin{cases} \frac{1}{y_F} \int_0^{y_F} \frac{x_F^\lambda}{x_I} dx_F & \text{if } x_I \geq y_F^\lambda, \\ \frac{1}{y_F} \int_0^{x_I^{1/\lambda}} \frac{x_F^\lambda}{x_I} dx_F + \frac{1}{y_F} \int_{x_I^{1/\lambda}}^{y_F} 1 \, dx_F & \text{if } x_I < y_F^\lambda. \end{cases} \\
&= \begin{cases} \frac{1}{y_F x_I}\left[\frac{1}{\lambda+1} x_F^{\lambda+1}\right]_0^{y_F} & \text{if } x_I \geq y_F^\lambda, \\ \frac{1}{\lambda+1} \frac{x_I^{\frac{\lambda+1}{\lambda}}}{y_F x_I} + 1 - \frac{x_I^{1/\lambda}}{y_F} & \text{if } x_I < y_F^\lambda. \end{cases} \\
&= \begin{cases} \frac{1}{\lambda+1} \frac{y_F^\lambda}{x_I} & \text{if } x_I \geq y_F^\lambda, \\ 1 - \frac{\lambda}{\lambda+1} \frac{x_I^{1/\lambda}}{y_F} & \text{if } x_I < y_F^\lambda. \end{cases}
\end{aligned}
$$

Then the fact that $r^*(x) = 1 - \min\{1, x_F^\lambda x_I^{-1}\}$ gives equation 2.2.

### A.1.6 Proof of lemma 5

*Proof.* Let $\mathcal{S} = (q, (w, f_e, f_0))$ be any scheme with finite support that deters bribe $b^*$. $\mathcal{S}$ and $b^*$ induce a distribution $p$ over payoffs defined by:

$$
p(v_F, v_I) := \mathbb{P}_q[v_F = b^* - r(x), v_I = f_e(x) - f_0(x) - b^*].
$$

The payoff distribution $p$ satisfies *Condition A* if $\max\{v_M, v_I\} \geq 0$ with probability 1, i.e. if at least one player benefits from the bribe ex post. Claim 7.1 shows that it is without loss of generality to restrict attention to such schemes.

**Claim 7.1.** *The government can restrict attention to schemes in which at least one player benefits from the bribe $b^*$ in each state: either $v_I = b^* - r(x) \geq 0$ or $v_F = f_e(x) - f_0(x) - b^* \geq 0$ for all messages $x \in X$.*

*Proof.* A proof by construction. Write $\Delta(x) := f_e(x) - f_0(x)$ and $\Delta'(x) := f_e'(x) - f_0'(x)$. Let $\mathcal{S}$ be a feasible scheme with message space $X = X_F \times X_I$, and an outcome $\tilde{x}$ in which $r(\tilde{x}) > b^*$ and $\Delta(\tilde{x}) < b^*$. Define a new scheme $\mathcal{S}'$ with

$$X_F' = X_F \cup \{x_F^a, x_F^b\} \setminus \{\tilde{x}_F\}$$

$$X_I' = X_I \cup \{x_I^a, x_I^b\} \setminus \{\tilde{x}_I\}$$

$$q'(x) = \begin{cases} q(x) & \text{if } x \in X \\ \frac{1}{2}q(\tilde{x}_i, x_{-i}) & \text{if } x \in \{x_i^a, x_i^b\} \times X_{-i} \text{ for some } i \\ \frac{1}{4}q(\tilde{x}_i) & \text{if } x \in \{x_F^a, x_F^b\} \times \{x_I^a, x_I^b\} \end{cases}$$

$$f_e'(x) = \begin{cases} f_e(x) & \text{if } x \in X \\ f_e(\tilde{x}_i, x_{-i}) & \text{if } \in \{x_i^a, x_i^b\} \times X_{-i} \text{ for some } i \\ f_e(\tilde{x}) - \frac{b - (f_e(\tilde{x}) - f_0(\tilde{x}))}{2} & \text{if } x \in \{(x_F^a, x_I^a), (x_F^b, x_I^b)\} \\ f_e(\tilde{x}) + \frac{b - (f_e(\tilde{x}) - f_0(\tilde{x}))}{2} & \text{if } x \in \{(x_F^a, x_I^b), (x_F^b, x_I^a)\} \end{cases}$$

$$f_0'(x) = \begin{cases} f_0(x) & \text{if } x \in X \\ f_0(\tilde{x}_i, x_{-i}) & \text{if } \in \{x_i^a, x_i^b\} \times X_{-i} \text{ for some } i \\ f_0(\tilde{x}) + \frac{b - (f_e(\tilde{x}) - f_0(\tilde{x}))}{2} & \text{if } x \in \{(x_F^a, x_I^a), (x_F^b, x_I^b)\} \\ f_0(\tilde{x}) - \frac{b - (f_e(\tilde{x}) - f_0(\tilde{x}))}{2} & \text{if } x \in \{(x_F^a, x_I^b), (x_F^b, x_I^a)\} \end{cases}$$

$$r'(x) = \begin{cases} r(x) & \text{if } x \notin \{x_F^a, x_F^b\} \times \{x_I^a, x_I^b\} \\ b^* & \text{if } x \in \{(x_F^a, x_I^a), (x_F^b, x_I^b)\} \\ 2r(\tilde{x}) - b^* & \text{if } x \in \{(x_F^a, x_I^b), (x_F^b, x_I^a)\}. \end{cases}$$

The new scheme satisfies limited liability because $r(\tilde{x}) > b^*$ implies $2r(x) - b^* > 0$ and $\Delta(\tilde{x}) < b^*$ implies $\Delta'_0(x) \leq b^* \leq \kappa$. The new scheme can be depicted as follows:

$$
\begin{array}{c|cccc}
\mathcal{S} & \cdots & \tilde{x}_I & \cdots \\
\hline
\vdots & \ddots & \vdots & \ddots \\
\tilde{x}_F & \cdots & (\Delta, r) & \cdots \\
\vdots & \ddots & \vdots & \ddots
\end{array}
\qquad \longmapsto \qquad
\begin{array}{c|ccccc}
\mathcal{S}' & \cdots & x_I^a & x_I^b & \cdots \\
\hline
\vdots & \ddots & \vdots & \ddots \\
x_F^a & \cdots & (2\Delta - b^*, b^*) & (b^*, 2r - b^*) & \cdots \\
x_F^b & \cdots & (b^*, 2r - b^*) & (2\Delta - b^*, b^*) & \cdots \\
\vdots & \ddots & \vdots & \ddots
\end{array}
$$

Suppose $\mathcal{S}'$ has an equilibrium $(\sigma'_F, \sigma'_I)$ for the bribe $b^*$. Consider the following cases:

1. If $\sigma'_F(x_F^a) = \sigma'_F(x_F^b)$ and $\sigma'_I(x_I^a) = \sigma'_I(x_I^b)$, then player $i$ get the same expected payoff from accepting the bribe in state $x_i^a$ as in state $x_i^b$, in scheme $\mathcal{S}$ for $i = F, I$. Define a strategy profile for the scheme $\mathcal{S}$ by $\sigma_i(x_i) = \sigma'_i(x_i)$ for all $x \neq \tilde{x}_i$, and $\sigma_i(\tilde{x}_i) = \sigma'_i(x_i^a)$ for both players $i = F, I$. Player $j$'s expected payoff from accepting the bribe when $i$ plays $\sigma_i$ in $\mathcal{S}$ is the same for all outcomes $x_j \neq \tilde{x}_j$ as when $i$ players $\sigma'_i$ in $\mathcal{S}'$. In outcome $\tilde{x}_j$, player $j$ gets the same expected payoff as in outcomes $x_j^a$ and $x_j^b$ in $\mathcal{S}'$. Thus, $\sigma_j$ must be a best response to $\sigma_i$ in $\mathcal{S}$ because $\sigma'_j$ is be a best response to $\sigma'_i$ in $\mathcal{S}'$. Therefore $(\sigma_F, \sigma_I)$ is an equilibrium under $\mathcal{S}$. The fact that $\mathcal{S}$ deters bribes implies that $\sigma_i = 0$ for either $i = I$ or $i = F$, which in turn implies that implies $\sigma'_i = 0$.

2. If $\sigma'_F(x_F^a) > \sigma'_F(x_F^b)$ then Ina gets a strictly higher expected payoff from reporting evidence in state $x_I^b$ than in state $x_I^a$ (because $2r(\tilde{x}) - b^* > b^*$, which implies that $\sigma'_I(x_I^a) > \sigma'_I(x_I^b)$. But then this implies that Finn has a higher expected return to suppressing evidence in state $x_F^b$ than in state $x_I^a$ (because $2f_e(\tilde{x}) - f_0(\tilde{x}) - b^* < b^*$, which means that $\sigma'_F(x_F^a) > \sigma'_F(x_F^b)$ cannot be a best response for Finn, giving a contradiction.

3. All other cases are analogous to case 2.

Therefore at least one agent plays the null strategy in every equilibrium of $\mathcal{S}'$. This means that $\mathcal{S}'$ is feasible. Moreover, $c(\mathcal{S}') = c(\mathcal{S})$, so $\mathcal{S}'$ costs the same as $\mathcal{S}$. $\qquad\square$

The *no-deal measure* of an equilibrium $(\sigma_F, \sigma_I)$ gives the probability that no bribery occurs. Formally, it is defined by $p_{\mathsf{ND}}(E) = \int_E 1 - \sigma_F(x_F)\sigma_I(x_I)\,dp(x)$ for all $E \subseteq \operatorname{supp} p$. A scheme $\mathcal{S}$ deters bribe $b^*$ if and only if the 'no-deal' measure of every equilibrium is equal to the prior $p$.

Given a prior distribution $p$ over payoffs, a distribution $\tilde{p}$ has a *disjoint negative decomposition bound* if there exists distributions $p_F$ and $p_I$ with the properties that

1. $\tilde{p} \leq p_F + p_I$,
2. $p_F + p_I \leq p$ and
3. $\int_{\mathbb{R}} v_F(x)\,dp_M(x) < 0$ and $\int_{\mathbb{R}} v_I(x)\,dp_I(x) < 0$.

**Proposition 9** (3.1, Carroll 2016). *If $p$ satisfies condition A then for every information structure there exists an equilibrium whose no-deal measure has a 'disjoint negative decomposition bound' (DNDB).*

The fact that $p$ deters bribes implies that $p$ has a DNDB.

**Proposition 10** (3.1, Carroll 2016). *If $p$ has a DNDB then there exists a public information structure such that, in any equilibrium, the no-deal measure is equal to $p$.*

Thus there exists a public information structure that deters bribe $b$. The cost of the scheme is unaffected by the information structure (it is determined exclusively by the payoff distribution), so the government can restrict attention to public schemes. $\qquad\square$

# Appendix B

# Inferring Inequality

## B.1 Mathematical Appendix

*Proof of lemma 6.* First we show that every distribution in either the median or the dominance boundaries (or both) is in the boundary. Let $\theta \in \bar{M} \cup \bar{D}$. Since $\theta = (F, G)$ is weakly ordered, $F$ and $G$ have at least one common median, $m$. The common median is unique unless $F_i = G_i = \frac{1}{2}$ for some $i \in [k]$. In this case we let $m = \min\{i \mid F_i = G_i = \frac{1}{2}\}$. Define a sequence $\{\theta^{1j}\}_{j \in \mathbb{N}}$ by $F^{1j} := \frac{1}{j}\mathcal{I}(j > m) + \left(1 - \frac{1}{j}\right)F$ and $G^{1j} := \frac{1}{j}\left(\frac{1}{2}\mathcal{I}(j > m) + \frac{1}{4}\right) + \left(1 - \frac{1}{j}\right)G$, for all $j < k$. This sequence converges to $\theta$ and it is easy to see graphically that each $\theta^{1j}$ is strictly ordered. Hence every element of $\bar{M} \cup \bar{D}$ is the limit of a sequence in the complement of the null space. Now define a sequence $\{\theta^{0j}\}_{j \in \mathbb{N}}$ by $F^{0j} := \frac{1}{j}\frac{1}{2} + \left(1 - \frac{1}{j}\right)F$ and $G^{0j} := \left(1 - \frac{1}{j}\right)G$.

This sequence also converges to $\theta$ and it is easy to see graphically that each $\theta^{1j}$ is strictly ordered, so long as there exists either an $i \leq m$ such that $F_i^0 = G_i^0$ or else an $i \geq m$ such that $F_i^0 > \frac{1}{2} = G_i^0$. If neither or these conditions hold then, in order for $\theta$ to be in $\bar{M} \cup \bar{D}$, there must exist either an $i > m$ such that $F_i^0 = G_i^0$ or else an $i \leq m$ such that $F_i^0 > \frac{1}{2} = G_i^0$. In this case we use the sequence defined by $G^{0j} := \frac{1}{j} + \left(1 - \frac{1}{j}\right)G$.

Now we show that every distribution in the boundary is in either the median or the dominance boundaries (or both). If $\theta \in (\bar{M} \cup \bar{D})^c$ is in neither the median nor dominance boundaries, then it must either be strictly ordered, or else unordered. Moreover, $\epsilon = \min\{|\frac{1}{2} - G_i|, |F_i - G_i| \mid i < k\}$ is strictly positive. If $\theta$ it strictly ordered then it strictly satisfies all the inequalities in definition 2 by a margin of at least $\epsilon$. Therefore any distribution $\theta'$ within a distance $\epsilon$ from $\theta$ must also strictly satisfy these inequalities. Thus there cannot exist any sequence of unordered distributions that converges to $\theta$. Similarly, if $\theta$ is unordered then it strictly violates at least one of the inequalities in definition 2 by a margin of at least $\epsilon$. Therefore any distribution $\theta'$ within a distance $\epsilon$ from $\theta$ must also violate the same inequality. Thus there cannot exist any sequence of ordered distributions that converges to $\theta$. $\square$

*Proof of lemma 7 point 3.* There are two ways the null hypothesis can be true: either one of the $k-1$ dominance conditions in [D1] or [D2] of definition 2 fails, or else the distributions do not share the same median and the conditions [M1] or [M2] in definition 2 fail. The easiest way for the former constraint to be satisfied is if $F_i = G_i$ for some $i \in [k-1]$ (which justifies a definition of strict MPS); the easiest way to satisfy the latter is if the median lies between two categories so that $G_{m-1} = \frac{1}{2}$ or $G_m = \frac{1}{2}$. Thus we can restate the problem:

$$\tilde{\theta} = \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y]$$

$$\text{s.t. } F_i = G_i \text{ for some } i \in [k-1]$$

$$\text{or } G_{m-1} = \frac{1}{2} \text{ or } G_m = \frac{1}{2}.$$

We now break the problem into two steps. We first find the $k+1$ distributions which maximise the likelihood, subject to each of these individual $k+1$ constraints, namely

$$\tilde{\theta}_i = \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y] \text{ s.t. } F_i = G_i \qquad\qquad \forall i < k \qquad\qquad \text{(B.1)}$$

$$\tilde{\theta}_k = \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y] \text{ s.t. } G_{m-1} = \frac{1}{2} \qquad\qquad\qquad \text{(B.2)}$$

$$\tilde{\theta}_{k+1} = \arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y] \text{ s.t. } G_m = \frac{1}{2} \qquad\qquad\qquad \text{(B.3)}$$

The solution to the original problem is then given by the distribution among these which maximises the sample's likelihood function: $\tilde{\theta} = \arg\max_{\tilde{\theta}=\tilde{\theta}_i} \mathbb{P}_{\tilde{\theta}_i}[x, y]$.

The solution to the problems in (B.1) are given in Davidson and Duclos (2013, p. 92) for each $i \in [k-1]$. The solution to (B.3) is found by noting that the independence of $f$ and $g$ implies that the solution to $\arg\max_{\theta \in \Theta_0} \mathbb{P}_\theta[x, y]$ s.t. $G_m = \frac{1}{2}$ is given by the pair $\left(\arg\max_f \mathbb{P}_f[x], \arg\max_g \mathbb{P}_g[y] \text{ s.t. } G_m = \frac{1}{2}\right)$. The first of these terms is simply. $x/n_x$. We solve for the second term by taking logarithms of the likelihood function $\mathbb{P}_g[y]$ (under the i.i.d. assumption) and by setting up the Lagrangian $\mathcal{L}(g, \lambda, \mu) = \sum_{i \in [k]} y_i \log g_i + \lambda(1 - \sum_{i \in [k]} g_i) + \mu(\frac{1}{2} - \sum_{i \in [m]} g_i)$. The first order condition requires that

$$\frac{\partial \mathcal{L}}{\partial g_i} = \begin{cases} \frac{y_i}{g_i} - \lambda - \mu & \text{if } i \leq m \\ \frac{y_i}{g_i} - \lambda & \text{if } i > m \end{cases} = 0$$

which implies

$$y_i = \begin{cases} (\lambda + \mu)\tilde{g}_i & \text{if } i \leq m \\ \lambda \tilde{g}_i & \text{if } i > m. \end{cases}$$

This in turn implies that $Y_m = \tilde{G}_m(\lambda + \mu) = \frac{1}{2}(\lambda + \mu)$ and $Y_k - Y_m = (1 - \tilde{G}_m)\lambda = \frac{1}{2}\lambda$. Together, these give $(\lambda + \mu) = 2Y_m$ and $\lambda = 2(Y_k - Y_m)$, and thus

$$
\tilde{g}_i = \begin{cases} \dfrac{y_i}{2Y_m} & \text{if } i \leq m \\ \dfrac{y_i}{2(n_y - Y_m)} & \text{if } i > m. \end{cases}
$$

The solution for (B.2) is found analogously. □

## B.2 Size and Power Curve Grid Points

Table B.1 lists the first coordinate of the DGP's on the median boundary and table B.2 lists the DGPs on the dominance boundary.[1] Our choice of dominance boundary DGPs depends on the sample size; the choice of median DGPs is the same for all sample sizes.

**Table B.1:** Median boundary DGPs used to construct the size curve in figure 3.3

| | Null DGPs on median boundary (all sample sizes) | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $f_1^{\text{med}}$ | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 | 0.5 |
| $g_1^{\text{med}}$ | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 |

**Table B.2:** Dominance boundary DGPs used to construct the size curve in figure 3.3

| | Null DGPs on dominance boundary (sample size dependent) | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $(n_x, n_y) = (10, 10), (100, 100), (1000, 1000)$ | | | | | | | | | | | |
| $f_1^{\text{dom}} = g_1^{\text{dom}}$ | | 0.46 | 0.43 | 0.39 | 0.36 | 0.32 | 0.28 | 0.25 | 0.21 | 0.16 | 0.10 | 0.05 |
| $(n_x, n_y) = (10, 100), (100, 1000)$ | | | | | | | | | | | |
| $f_1^{\text{dom}} = g_1^{\text{dom}}$ | 0.47 | 0.45 | 0.42 | 0.39 | 0.36 | 0.34 | 0.30 | 0.27 | 0.23 | 0.15 | 0.10 | 0.05 |
| $(n_x, n_y) = (10, 1000)$ | | | | | | | | | | | |
| $f_1^{\text{dom}} = g_1^{\text{dom}}$ | 0.48 | 0.45 | 0.43 | 0.40 | 0.37 | 0.35 | 0.32 | 0.29 | 0.25 | 0.15 | 0.10 | 0.05 |
| $(n_x, n_y) = (100, 10), (1000, 100)$ | | | | | | | | | | | |
| $f_1^{\text{dom}} = g_1^{\text{dom}}$ | | 0.45 | | 0.40 | 0.36 | 0.31 | 0.26 | 0.21 | 0.17 | 0.12 | 0.07 | 0.05 |
| $(n_x, n_y) = (1000, 10)$ | | | | | | | | | | | |
| $f_1^{\text{dom}} = g_1^{\text{dom}}$ | | 0.45 | | 0.40 | 0.35 | | 0.30 | 0.25 | 0.20 | 0.15 | 0.10 | 0.05 |

Table B.3 lists the first coordinate of the DGPs on the interior locus of the alternative space used to evaluate power. The concentrated distributions $f$ are chosen to be the same for all sample sizes $n_x : n_y$; the precise choice of spread distribution $g$ is then chosen to ensure that it is equidistant from the median and dominance boundaries.

---

1. Because $k = 2$, it is sufficient to note only the mass in the first category; the full distributions can be recovered from $f_1$ and $g_1$ by $f = (f_1, 1 - f_1)$ and $g = (g_1, 1 - g_1)$.

Table B.3: Alternative DGPs on the interior locus illustrated in figure 3.4.

| | Alternative DGPs (spread distributions depends on sample size) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $f_1^A$ | 0.05 | 0.1 | 0.15 | 0.2 | 0.25 | 0.3 | 0.35 | 0.4 | 0.45 |
| (nx,ny)=(10,10),(100,100),(1000,1000) | | | | | | | | | |
| $g_1^A$ | 0.277 | 0.312 | 0.342 | 0.368 | 0.392 | 0.415 | 0.437 | 0.458 | 0.479 |
| (nx,ny)=(10,100),(100,1000) | | | | | | | | | |
| $g_1^A$ | 0.252 | 0.289 | 0.32 | 0.349 | 0.376 | 0.402 | 0.427 | 0.452 | 0.476 |
| (nx,ny)=(10,100),(10,1000) | | | | | | | | | |
| $g_1^A$ | 0.244 | 0.282 | 0.314 | 0.344 | 0.371 | 0.398 | 0.424 | 0.450 | 0.475 |
| (nx,ny)=(100,10) | | | | | | | | | |
| $g_1^A$ | 0.235 | 0.280 | 0.315 | 0.346 | 0.375 | 0.401 | 0.427 | 0.451 | 0.476 |
| (nx,ny)=(1000,10) | | | | | | | | | |
| $g_1^A$ | 0.219 | 0.268 | 0.306 | 0.339 | 0.369 | 0.397 | 0.424 | 0.449 | 0.475 |

# Bibliography

Abul Naga, R., & Stapenhurst, C. (2015). Estimation of inequality indices of the cumulative distribution function. *Economics Letters*, *130*, 109-112.

Abul Naga, R., Stapenhurst, C., & Yalonetzky, G. (2020). Asymptotic versus bootstrap inference for inequality indices of the cumulative distribution function. *Econometrics*, *8*(1), 8.

Abul Naga, R., & Yalcin, T. (2008). Inequality measurement for ordered response health data. *Journal of Health Economics*, *27*, 1614-25.

Akerlof, G. A. (1970). The market for "lemons": Quality uncertainty and the market mechanism. *The Quarterly Journal of Economics*, *84*(3), 488–500. Retrieved from `http://www.jstor.org/stable/1879431`

Allison, R. A., & Foster, J. E. (2004). Measuring health inequality using qualitative data. *Journal of Health Economics*, *23*, 505-24.

Aoyagi, M. (2005). Collusion through mediated communication in repeated games with imperfect private monitoring. *Economic Theory*, *25*(2), 455–475. Retrieved from `http://www.jstor.org/stable/25055890`

Apouey, B. (2007). Measuring health polarisation with self-assessed health data. *Health Economics*, *16*, 875-94.

Balestra, C., & Ruiz, N. (2015). Scale-invariant measurement of inequality and welfare in ordinal achievements: an application to subjective well-being and education in oecd countries. *Social Indicators Research*, *123*(2), 479–500.

Baliga, S., & Sjöström, T. (1998). Decentralization and collusion. *Journal of Economic Theory*, *83*(2), 196-232. Retrieved from `https://www.sciencedirect.com/science/article/pii/S002205319692462X` doi: https://doi.org/10.1006/jeth.1996.2462

Beck, J. (2000). The false claims act and the english eradication of qui tam legislation. *North Carolina Law Review*, *78*(3), 539–642.

Becker, G. S. (1968). Crime and punishment: An economic approach. *Journal of Political Economy*, *76*(2), 169-217. Retrieved from `https://doi.org/10.1086/259394` doi: 10.1086/259394

Ben-Porath, E., & Kahneman, M. (1996). Communication in repeated games with private monitoring. *Journal of Economic Theory*, *70*(2), 281-297. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0022053196900903` doi: https://doi.org/10.1006/jeth.1996.0090

Ben-Porath, E., & Lipman, B. (2012). Implementation with partial provability. *Journal of Economic Theory*, *147*(5), 1689-1724. Retrieved from `https://EconPapers.repec.org/RePEc:eee:jetheo:v:147:y:2012:i:5:p:1689-1724`

Bergemann, D., & Morris, S. (2019, March). Information design: A unified perspective. *Journal of Economic Literature*, *57*(1), 44-95. Retrieved from `https://www.aeaweb.org/articles?id=10.1257/jel.20181489` doi: 10.1257/jel.20181489

Berger, R. (1982). Multiparameter hypothesis testing and acceptance sampling. *Technometrics*, *24*, 295-300.

Carroll, G. (2016). Informationally robust trade and limits to contagion. *Journal of Economic Theory*, *166*(C), 334-361. Retrieved from `https://EconPapers.repec.org/RePEc:eee:jetheo:v:166:y:2016:i:c:p:334-361`

Chakravarty, S., & Maharaj, B. (2012, May). Ethnic polarization orderings and indices. *Journal of Economic Interaction and Coordination*, *7*(1), 99-123. Retrieved from `https://ideas.repec.org/a/spr/jeicoo/v7y2012i1p99-123.html` doi: 10.1007/s11403-011-0084-z

Chakravarty, S. R., & Maharaj, B. (2015). Generalized gini polarization indices for an ordinal dimension of human well-being. *International Journal of Economic Theory*, *11*(2), 231-246. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/ijet.12062` doi: https://doi.org/10.1111/ijet.12062

Che, Y.-K., & Kim, J. (2006). Robustly collusion-proof implementation. *Econometrica*, *74*(4), 1063-1107. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1468-0262.2006.00694.x` doi: 10.1111/j.1468-0262.2006.00694.x

Cremer, J., & McLean, R. P. (1988, November). Full Extraction of the Surplus in Bayesian and Dominant Strategy Auctions. *Econometrica*, *56*(6), 1247-1257. Retrieved from `https://ideas.repec.org/a/ecm/emetrp/v56y1988i6p1247-57.html`

Crémer, J. (1996). Manipulations by coalitions under asymmetric information: The case of groves mechanisms. *Games and Economic Behavior*, *13*(1), 39-73. Retrieved from `https://www.sciencedirect.com/science/article/pii/S089982569690024X` doi: https://doi.org/10.1006/game.1996.0024

Davidson, R., & Duclos, J.-Y. (2013). Testing for restricted stochastic dominance. *Econometric Reviews*, *32*(1), 84-125. Retrieved from `https://doi.org/10.1080/07474938.2012.690332` doi: 10.1080/07474938.2012.690332

Davidson, R., & MacKinnon, J. G. (1998). Graphical methods for investigating the size and power of hypothesis tests. *The Manchester School*, *66*(1), 1-26. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/1467-9957.00086` doi: https://doi.org/10.1111/1467-9957.00086

Davison, A., & Hinkley, D. (1997). *Bootstrap methods and their applications*. Cambridge: Cambridge University Press.

Duflo, E., Greenstone, M., Pande, R., & Ryan, N. (2013, 09). Truth-telling by Third-party Auditors and the Response of Polluting Firms: Experimental Evidence from India*. *The Quarterly Journal of Economics*, *128*(4), 1499-1545. Retrieved from `https://doi.org/10.1093/qje/qjt024` doi: 10.1093/qje/qjt024

Dutta, I., & Foster, J. (2013). Inequality of happiness in the u.s.: 1972–2010. *The Review of Income and Wealth*, *59*(3), 393-415.

Felli, L., & Hortala-Vallve, R. (2016, 11). Collusion, blackmail and whistle-blowing *. *Quarterly Journal of Political Science*, *11*. doi: 10.1561/100.00015060

Forge, F., & Serrano, R. (2013). Cooperative games with incomplete information: Some open problems. *International Game Theory Review*, *15*(02), 1340009. Retrieved from `https://doi.org/10.1142/S0219198913400094` doi: 10.1142/S0219198913400094

Garrett, D. F., Georgiadis, G., Smolin, A., & Szentes, B. (2020). Optimal technology design. *Available at SSRN 3720594*.

Green, J., & Laffont, J.-J. (1979, 04). On Coalition Incentive Compatibility. *The Review of Economic Studies*, *46*(2), 243-254. Retrieved from `https://doi.org/10.2307/2297048` doi: 10.2307/2297048

Green, J. R., & Laffont, J.-J. (1986). Partially verifiable information and mechanism design. *The Review of Economic Studies*, *53*(3), 447–456.

Gunawan, D., Griffiths, W. E., & Chotikapanich, D. (2018). Bayesian inference for health inequality and welfare using qualitative data. *Economics Letters*, *162*, 76-80. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0165176517304603` doi: https://doi.org/10.1016/j.econlet.2017.11.005

Halac, M., Lipnowski, E., & Rappoport, D. (2021, March). Rank uncertainty in organizations. *American Economic Review*, *111*(3), 757-86. Retrieved from `https://www.aeaweb.org/articles?id=10.1257/aer.20200555` doi: 10.1257/aer.20200555

Holmström, B. (1979). Moral hazard and observability. *The Bell Journal of Economics*, *10*(1), 74–91. Retrieved from `http://www.jstor.org/stable/3003320`

Kajii, A., & Morris, S. (1997). The robustness of equilibria to incomplete information. *Econometrica*, *65*(6), 1283–1309. Retrieved from `http://www.jstor.org/stable/2171737`

Kobus, M. (2015). Polarisation measurement for ordinal data. *Journal of Economic Inequality*, *13*(2), 275-97.

Kobus, M., & Milos, P. (2012). Inequality decomposition by population subgroups for ordinal data. *Journal of Health Economics*, *31*, 15-21.

Koessler, F., & Perez-Richet, E. (2019). Evidence reading mechanisms. *Social Choice and Welfare*, *53*(3), 375–397.

Kofman, F., & Lawarrée, J. (1993). Collusion in hierarchical agency. *Econometrica*, *61*(3), 629–656. Retrieved from `http://www.jstor.org/stable/2951721`

Laffont, J.-J., & Martimort, D. (1997). Collusion under asymmetric information. *Econometrica*, *65*(4), 875-912. Retrieved from `https://EconPapers.repec.org/RePEc:ecm:emetrp:v:65:y:1997:i:4:p:875-912`

Laffont, J.-J., & Martimort, D. (2000). Mechanism design with collusion and correlation. *Econometrica*, *68*(2), 309-342. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/1468-0262.00111` doi: 10.1111/1468-0262.00111

Langpap, C., & Shimshack, J. P. (2010). Private citizen suits and public enforcement: Substitutes or complements? *Journal of Environmental Economics and Management*, *59*, 235–249.

Latham, K., & Peek, C. (2013). Self-rated health and morbidity onset among late midlife u.s. adults. *Journals of Gerontology Series B: Psychological Sciences and Social Sciences*, *68*(1), 107-116.

Lazar, A., & Silber, J. (2013). On the cardinal measurement of health inequality when only ordinal information is available on individual health status. *Health Economics*, *22*(1), 106-13.

Lehmann, E., & Romano, J. (2005). *Testing statistical hypotheses*. Springer.

Liu, Z. J., Wang, Z., & Yin, Z. (2022). When is duplication of effort a good thing in law enforcement? *Journal of Public Economic Theory*, *n/a*(n/a). Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/jpet.12568` doi: https://doi.org/10.1111/jpet.12568

Madden, D. (2010). Ordinal and cardinal measures of health inequality: An empirical comparison. *Health Economics*, *19*(2), 243-250.

Mathevet, L., Perego, J., & Taneva, I. (2020). On information design in games. *Journal of Political Economy*, *128*(4), 1370-1404. Retrieved from `https://doi.org/10.1086/705332` doi: 10.1086/705332

McAfee, R. P., Mialon, H. M., & Mialon, S. H. (2008). Private v. public antitrust enforcement: A strategic analysis. *Journal of Public Economics*, *92*(10), 1863 - 1875. Retrieved from `http://www.sciencedirect.com/science/article/pii/S0047272708000637` doi: https://doi.org/10.1016/j.jpubeco.2008.04.005

Mendelson, H. (1987). Quantile- preserving spread. *Journal of Economic Theory*, *42*, 334-51.

Mood, A., Graybill, F., & Boes, D. (1974). *Introduction to the theory of statistics*. McGraw-Hill.

Morris, S., & Ui, T. (2005). Generalized potentials and robust sets of equilibria. *Journal of Economic Theory*, *124*(1), 45-78. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0022053104002182` doi: https://doi.org/10.1016/j.jet.2004.06.009

Myerson, R. B. (2007). Virtual utility and the core for games with incomplete information. *Journal of Economic Theory*, *136*(1), 260-285. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0022053106001438` doi: https://doi.org/10.1016/j.jet.2006.08.002

Ortner, J., & Chassang, S. (2018). Making corruption harder: Asymmetric information, collusion, and crime. *Journal of Political Economy*, *126*(5), 2108-2133. Retrieved from `https://doi.org/10.1086/699188` doi: 10.1086/699188

Polinsky, A. M. (1980). Private versus public enforcement of fines. *The Journal of Legal Studies*, *9*(1), 105–127. Retrieved from `http://www.jstor.org/stable/724040`

Rahman, D. (2012, May). But who will monitor the monitor? *American Economic Review*, *102*(6), 2767-97. Retrieved from `http://www.aeaweb.org/articles?id=10.1257/aer.102.6.2767` doi: 10.1257/aer.102.6.2767

Reardon, S. (2009). Measures of ordinal segregation. In Y. Fluckiger, S. Reardon, & J. Silber (Eds.), *Occupational and residential segregation* (Vol. 17).

Rothschild, M., & Stiglitz, J. (1970). Increasing risk: I. a definition. *Journal of Economic Theory*, *2*(3), 225-43.

Seth, S., & Yalonetzky, G. (2021). Assessing deprivation with an ordinal variable: Theory and application to sanitation deprivation in bangladesh. *The World Bank Economic Review*, *35*(3), 793–811.

Silber, J., & Yalonetzky, G. (2021, July). Measuring welfare, inequality and poverty with ordinal variables. In K. Zimmermann (Ed.), *Handbook of labor, human resources and population economics.* Springer. Retrieved from `https://eprints.whiterose.ac.uk/171816/` (This item is protected by copyright. This is an author produced version of a book chapter published in Handbook of Labor, Human Resources and Population Economics. Uploaded in accordance with the publisher's self-archiving policy.)

Stapenhurst, C. (2019). *How Many Corruptible Monitors does it take to Implement an Action?* (Unpublished doctoral dissertation). University of Edinburgh.

Stevens, S. S. (1946). On the theory of scales of measurement. *Science, New Series*, *103*(2684), 677–680.

Strausz, R. (1997). Delegation of monitoring in a principal-agent relationship. *Review of Economic Studies*, *64*(3), 337-357. Retrieved from `https://EconPapers.repec.org/RePEc:oup:restud:v:64:y:1997:i:3:p:337-357`.

Taneva, I. (2019, November). Information design. *American Economic Journal: Microeconomics*, *11*(4), 151-85. Retrieved from `https://www.aeaweb.org/articles?id=10.1257/mic.20170351` doi: 10.1257/mic.20170351

Tirole, J. (1986). Hierarchies and bureaucracies: On the role of collusion in organizations. *Journal of Law, Economics, & Organization*, *2*(2), 181–214. Retrieved from `http://www.jstor.org/stable/765048`

Topsøe, F. (2004). Some bounds for the logarithmic function. *RGMIA Res. Rep. Collection*, *7*(2), 1–20.

United Nations. (2018). *Global cost of corruption at least 5 per cent of world gross domestic product, secretary-general tells security council, citing world economic forum data.* Retrieved from `https://www.un.org/press/en/2018/sc13493.doc.htm`

Vafaï, K. (2005). Collusion and organization design. *Economica*, *72*(285), 17-37. Retrieved from `https://onlinelibrary.wiley.com/doi/abs/10.1111/j.0013-0427.2005.00400.x` doi: 10.1111/j.0013-0427.2005.00400.x

Vanstraelen, A., Richard, C., & R. Francis, J. (2009, 11). Assessing france's joint audit requirement: Are two heads better than one? *Auditing A Journal of Practice & Theory*, *28*. doi: 10.2308/aud.2009.28.2.35

von Negenborn, C., & Pollrich, M. (2020). Sweet lemons: Mitigating collusion in organizations. *Journal of Economic Theory*, *189*, 105074. Retrieved from `https://www.sciencedirect.com/science/article/pii/S0022053120300703` doi: https://doi.org/10.1016/j.jet.2020.105074

Yalonetzky, G. (2013). Stochastic dominance with ordinal variables: Conditions and a test. *Econometric Reviews*, *32*(1), 126-63.