

Why do people punish? Evidence for a range of strategic concerns

Scott Claessens^{*1}, Quentin D. Atkinson¹, & Nichola Raihani^{1,2}

¹ School of Psychology, University of Auckland, Auckland, New Zealand

² Department of Experimental Psychology, University College London, London, United Kingdom

* Correspondence concerning this article should be addressed to Scott Claessens, Level 2, Building 302, 23 Symonds Street, Auckland, New Zealand. E-mail: scott.claessens@gmail.com

This manuscript has not yet been peer-reviewed.

This manuscript was published as a pre-print on PsyArXiv under a CC BY 4.0 license:

<https://doi.org/10.31234/osf.io/ys6rm>

Classification: Social Sciences; Psychological and Cognitive Sciences

Keywords: punishment; cooperation; economic games

Abstract

Costly punishment is thought to be a key mechanism sustaining human cooperation. However, the motives for punitive behaviour remain unclear. Although often assumed to be motivated by a desire to convert cheats into cooperators, punishment is also consistent with other functions, such as levelling payoffs or improving one's relative position. We used six economic games to tease apart different motives for punishment and to explore whether different punishment strategies were associated with personality variables, political ideology, and religiosity. We used representative samples from the United Kingdom and the United States ($N = 2010$) to estimate the frequency of different punishment strategies in the population. The most common strategy was to never punish: prosocial individuals and those scoring high in honesty-humility were among the least punitive in our sample, while religious individuals and those scoring high for right-wing authoritarianism and social dominance orientation were more punitive. For people who did punish, strategy use was more consistent with egalitarian motives than behaviour-change motives. Nevertheless, different punishment strategies were also associated with personality, social preferences, political ideology, and religiosity. Self-reports of behaviour in the games suggested that people have some insight into their punishment strategy. These findings highlight the multipurpose nature of human punishment and show how the different motives underpinning punishment decisions are linked with core character traits.

Introduction

Humans cooperate on a scale that is unparalleled in the animal kingdom. One mechanism thought to sustain this level of cooperation is costly punishment, whereby individuals harm others at a personal cost (1), ostensibly encouraging cooperative behaviour from the target, or bystanders (2–4), in the future. Punishment offers a route to maintaining or increasing cooperation by changing the payoff structure of social interactions such that it no longer pays to cheat or exploit social partners (1, 5). Yet, despite its theoretical importance, the question of why people choose to punish others is still hotly contested (6). In this study, we use a battery of economic games to disentangle the different motives underpinning punishment and explore how these motives vary across individuals.

In humans, many studies of punishment have been carried out in laboratory settings using economic games (6–15). In these games, participants are usually given a sum of money that they can use to invest in collective action or to help others. Alternatively, participants can ‘cheat’ by keeping the money for themselves or by exploiting the contributions of others. Punishment is introduced into such games by giving participants the option to pay a small ‘fee’ to impose a greater ‘fine’ on their co-players. Several lines of experimental evidence indicate that people use this punishment option (13), that they enjoy punishing (16), and that they frequently, though not always (17), punish cheating or exploitative co-players (11, 12).

Evidence from these experiments suggests that the threat of costly punishment plays an important role in promoting human cooperation. People tend to cooperate more in games where punishment is possible compared to those where it is not (6–8). The effect that the threat of punishment has on cooperation is also evident in the higher contributions typically observed in the Ultimatum Game (where punishment is possible) compared to the structurally-similar Dictator Game (where it is not) (18). This typical

cooperation-enhancing effect of punishment has also been observed across societies (8), leading some to suggest that costly punishment has played a key role in the cultural evolution of cooperation in humans (19–22).

Nevertheless, it remains unclear whether individuals playing economic games use punishment as a behaviour-change tool to enforce cooperation or as a means to achieve other ends. Some have argued that punishment is primarily used to shape future behaviour, either to deter personal harm (3, 10, 23) or to uphold normative standards of cooperative behaviour (20, 21, 24–27). But while the *threat* of punishment can have a cooperation-enhancing effect, the *enactment* of this punishment does not consistently deter targets from cheating in the future (6). This calls into question whether punishment primarily operates as a behaviour-change tool or whether it is used to achieve other goals.

Beyond behaviour shaping concerns, there are a host of other reasons why people may want to punish in economic games. Punishers might be motivated by a desire for retribution rather than deterrence, punishing in proportion to the amount of harm that was personally caused (28). Punishment might be driven by concerns about relative payoffs, such as disadvantageous inequity aversion (i.e., avoiding having less than others) (6, 29) and/or general egalitarian preferences (i.e., wanting all participants to receive the same payoffs) (30). Such concerns about relative payoffs may be activated when participants earn less than cheats in economic games or when there are income disparities in these settings. People might also use punishment for competitive purposes, seeking advantageous inequity for themselves (i.e., having more than others) and/or improving their relative position (6).

Common economic game designs have been unable to tease apart these different motives for punishment because participants who interact with cheats in these games experience both losses *and* lower relative payoffs. The typical 1:3 fee-fine ratio of punishment in economic games compounds this issue. With this setup, people can simultaneously use punishment to reciprocate losses, to deter others from cheating, and to

reduce or reverse disparities in payoffs between themselves and targets. To add to this complexity, it is evident that people use punishment in seemingly disparate ways: punishing when no behaviour change is possible, such as in one-shot games (13, 29, 31, 32), on the very last round of repeated games (33), or in games where the target never learns about the punishment (34); punishing those who did not cheat or who over-contributed to collective action (antisocial punishment) (17, 35); punishing in scenarios where they were not personally harmed (third-party punishment) (36); and punishing in scenarios where disparities in payoffs did not arise from participants' actions (30, 37, 38).

The general conclusion from this research is that there is no one unifying function of costly punishment in humans. Instead, punishment should be thought of as a flexible behavioural tool that serves a variety of functions that are not mutually exclusive (6). Due to its multipurpose nature, we should therefore expect variation in punishment strategies in the population, much like the observed variation in social learning strategies (39). Some individuals may use punishment as a behaviour shaping tool, for example, while others may use it to reduce or reverse payoff differentials.

This insight raises several underexplored questions. First, which punishment strategies are more frequent in human populations? Second, what traits predict adherence to a particular strategy? Previous work has reported that personality is related to cooperative behaviour (40) and demographics, political ideology, and religiosity are related to punitive behaviour (41), but no research has related these variables to specific punishment strategies. Third, do people have insight into their own punishment strategy? Previous work has argued that people are often unaware of the underlying function of their punitive behaviour, yet they feel compelled to enact it anyway (28, 42).

Here, we aim to delineate nine possible punishment strategies by asking whether people punish in a manner consistent with a specific strategy and, if so, what other characteristics (personality, social preferences, political orientation, religious views) predict

the use of different punishment strategies. Table 1 summarises the potential functions for costly punishment in the economic games that we considered, and the behavioural strategies they predict. Note that Table 1 is not an exhaustive list of all possible punishment strategies: we do not include reputational functions of punishment in this table, such as signalling trustworthiness (4, 43–46), because our focus is on punishment strategies in anonymous economic games without reputational incentives (but see ref(47)).

Building on previous designs (29, 31, 48, 49), we employ a suite of one-shot economic games where individuals are given the opportunity to punish targets at a personal cost (Figure 1). In each game, targets either steal from another individual or do nothing. Representative samples of participants from the United Kingdom ($n = 1014$) and the United States ($n = 996$) completed all six games on the online platform Prolific. We carefully designed the suite of games to tease apart the proposed punishment strategies in Table 1, such that each strategy predicts a different pattern of behaviour across all the games (see Methods for more detail about the six games). We use the resulting behavioural patterns to discern which punishment strategy participants are employing. We then combine these behavioural patterns with data on demographics, personality, social preferences, political ideology, religiosity, and self-reported strategy usage.

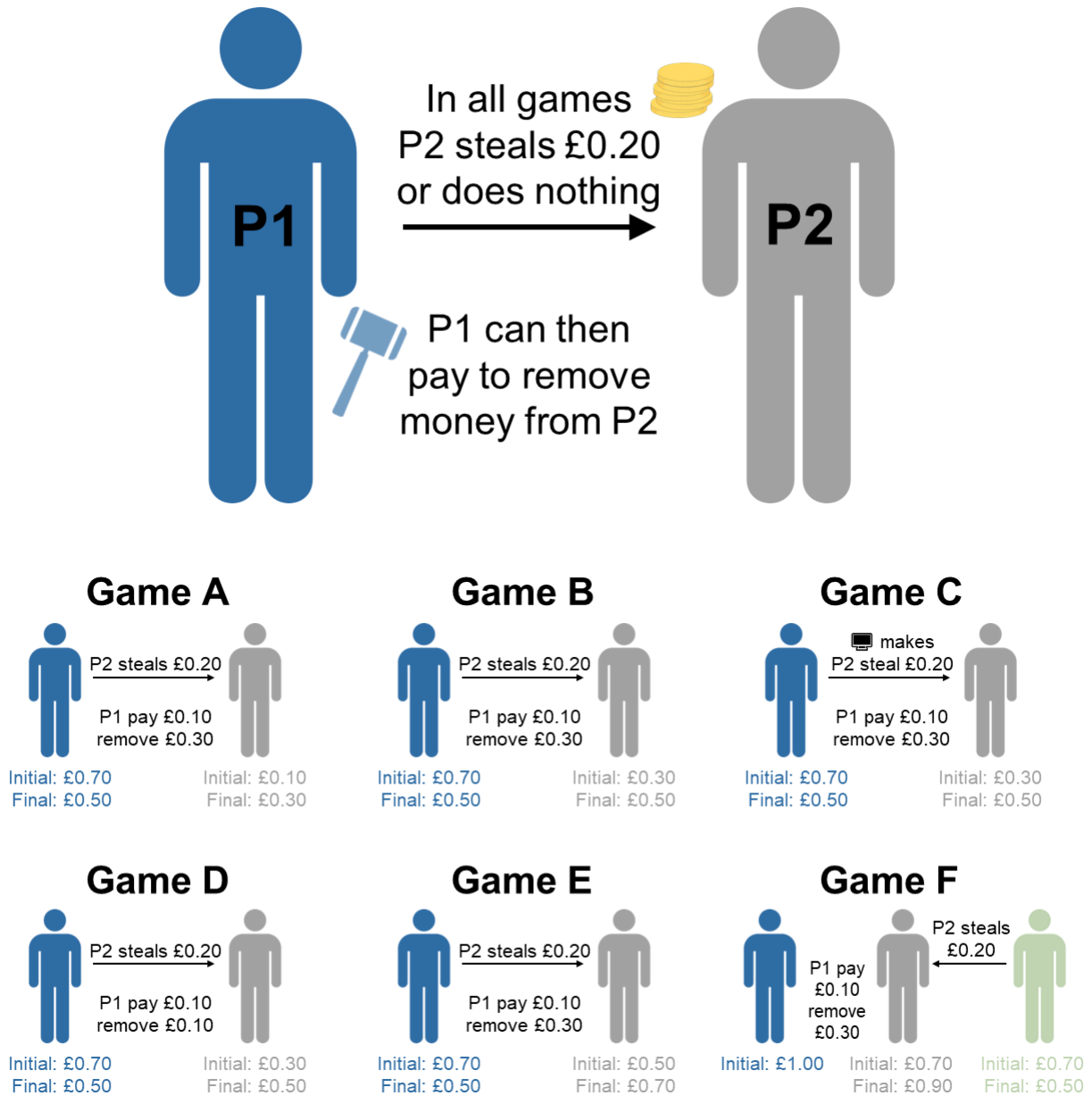


Figure 1. Visual summary of the six economic games. In the games, Player 2 either steals £0.20 from Player 1 (the focal player) or does nothing. Player 1 is then given the option to punish by paying a certain amount of money to remove money from Player 2 (this money is destroyed). The six games are variants on this general setup, creating situations where (A) Player 2 is still worse off after stealing, (B) Player 2 creates equality by stealing, (C) the computer “decides” whether Player 2 steals, (D) the fee-fine ratio is 1:1, (E) Player 2 is better off after stealing, and (F) Player 2 steals instead from a third-party.

Table 1

Summary of the different functions for punishment and the behavioural strategies they predict. Games A-F are the games employed in the current study (see Methods for more details). In each of the six games, participants are given the opportunity to punish players who “steal” and those who do not, meaning that participants make twelve punishment decisions in total. Each behavioural strategy implies a unique pattern of punishment across all decisions (see Methods for detailed explanation of strategies). Green ticks reflect decisions to punish, red crosses reflect decisions to not punish. In column headers, payoffs at the first stage (above) and the second stage (below) are denoted as P1-P2 (or P2-P3 [P1] for Game F) where participants take the role of P1 and P2 is the target of punishment. AI = advantageous inequity, DI = disadvantageous inequity.

Function	Behavioural strategy	Game A (AI) 70-10			Game B (Equal) 70-30			Game C (Computer) 70-30			Game D (1:1 Fee-Fine) 70-30			Game E (DI) 70-50			Game F (Third-Party) 70-70 [100]		
		Steal 50-30	No steal 70-10		Steal 50-50	No steal 70-30		Steal 50-50	No steal 70-30		Steal 50-50	No steal 70-30		Steal 50-70	No steal 70-50		Steal 50-90 [100]	No steal 70-70 [100]	
Deterrent	Punish to deter another who has harmed you from harming you again in the future	✓	✗		✓	✗		✗	✗		✓	✗		✓	✗		✗	✗	
Norm-enforcing	Punish to enforce a shared anti-harm norm and encourage future norm compliance, even amongst third parties	✓	✗		✓	✗		✗	✗		✓	✗		✓	✗		✓	✗	
Retributive	Punish if doing so harms another who has harmed you	✓	✗		✓	✗		✓	✗		✓	✗		✓	✗		✗	✗	
Avoid DI	Punish if doing so avoids disadvantageous inequity for self	✗	✗		✗	✗		✗	✗		✗	✗		✓	✗		✗	✗	
Egalitarian	Punish if doing so makes payoffs for all more equal	✗	✗		✗	✗		✗	✗		✗	✗		✓	✗		✓	✗	
Seek AI	Punish if doing so produces advantageous inequity for self	✗	✗		✓	✗		✓	✗		✗	✗		✗	✗		✗	✗	
Competitive	Punish if doing so improves your relative position	✓	✓		✓	✓		✓	✓		✗	✗		✓	✓		✓	✓	
Antisocial	Punish exclusively those who do not cause harm	✗	✓		✗	✓		✓	✓		✗	✓		✗	✓		✗	✓	
Never punish	Never punish others	✗	✗		✗	✗		✗	✗		✗	✗		✗	✗		✗	✗	

Results

The overall pattern of punitive behaviour in the six economic games was very similar across both countries (Figure 2). Participants were generally more likely to punish targets who stole compared to targets who did not steal (multilevel logistic regression; $b = 1.93$, standard error = 0.27, $p < .001$). Participants were also more likely to punish when targets' stealing behaviour generated inequalities, specifically in Games E and F ($b = 2.42$, SE = 0.44, $p < .001$).

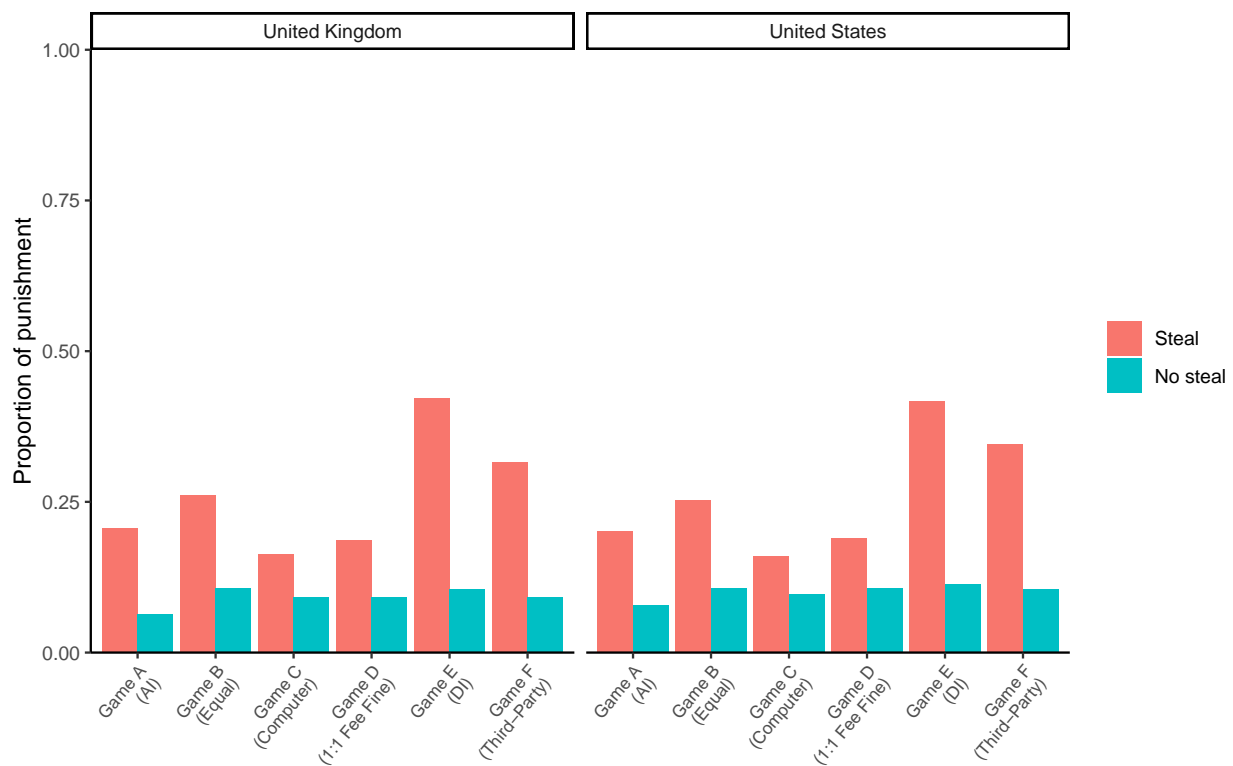


Figure 2. Overall pattern of punitive behaviour across all six economic games, split by country. AI = advantageous inequity, DI = disadvantageous inequity.

We classified participants into a particular strategy if their behaviour across all twelve decisions matched our behavioural predictions shown in Table 1 exactly. Table 2 shows the proportion of participants following each strategy, with N/A used to represent participants who did not fit exactly into any particular strategy type. Overall, 59% of our participants could be classified exactly into one of the strategies. The most common

Table 2

Counts and proportions of participants following each punishment strategy exactly, split by country.

N/A implies that participants were unable to be classified exactly into any of the punishment strategies. AI = advantageous inequity, DI = disadvantageous inequity.

Strategy	United Kingdom (N = 1014)		United States (N = 996)	
	N	Prop	N	Prop
Deterrent	9	0.009	6	0.006
Norm-enforcing	8	0.008	16	0.016
Retributive	6	0.006	5	0.005
Avoid DI	67	0.066	62	0.062
Egalitarian	65	0.064	71	0.071
Seek AI	2	0.002	0	0.000
Competitive	3	0.003	1	0.001
Antisocial	0	0.000	0	0.000
Never punish	426	0.420	447	0.449
N/A	428	0.422	388	0.390

strategy in both countries was to never punish across any of the games. The next most common strategies were those that care about minimising payoff differences (avoid disadvantageous inequity, egalitarian). Less common were the behaviour shaping strategies (deterrent, norm-enforcing), the retributive strategy, and the competitive strategies (seek advantageous inequity, competitive). Although participants often punished targets who did not steal in the six games (Figure 2), no participants followed the antisocial strategy by exclusively punishing targets who did not steal across *all* games.

To further investigate the strategies that participants were following, we examined the most common patterns of punitive behaviour across all twelve decisions. Supplementary Table S1 shows the proportion of participants following the 25 most common behavioural patterns, including, where appropriate, the predetermined strategies from Table 1. In both countries, a common pattern of behaviour not captured by any of

the strategies was punishing only when the target stole in the third-party game (Game F). Punishment in this game is consistent with an egalitarian motive, as stealing produces unequal outcomes, but third-party punishment here is also consistent with norm-enforcing and competitive motives (see Table 1). Other common behavioural patterns not captured by our strategies included punishing whenever the target stole across all games and always punishing in every game irrespective of the targets' behaviour.

While it is useful to look at exact patterns of behaviour, participants may not have implemented their chosen punishment strategy with exact precision. In reality, strategies may have been implemented probabilistically for each punishment decision. There is also the possibility of implementation errors, whereby participants occasionally “slip up” and make decisions that are incongruent with a particular strategy. This may explain why some participants were unable to be classified exactly into a single punishment strategy.

To deal with this complexity and include all observed data in our frequency estimates, we fitted a Bayesian latent state model to the data. This model assumes that the nine strategies in Table 1 (plus a “random choice” strategy that chooses randomly for each decision) are the only latent strategies and that these are instantiated into observed behaviour according to the logic in Table 1 with some probability of implementation error (i.e., an intention to punish is implemented as non-punishment and vice versa). Averaging over all strategies and incorporating the possibility of implementation errors, the model estimates the probability of participants following any particular strategy, given the observed data.

The posterior estimates from the model are presented in Figure 3. The posterior probabilities for each strategy did not differ between the two countries. In both countries, the never punish strategy had the highest probability, followed by the egalitarian strategy. The norm-enforcing and seek advantageous inequity strategies were the next most likely, with higher posterior estimates than the competitive and antisocial strategies. None of the

other strategies differed in their posterior estimates. The same general pattern emerged when we analysed the full dataset without pre-registered exclusions (Supplementary Figure S1).

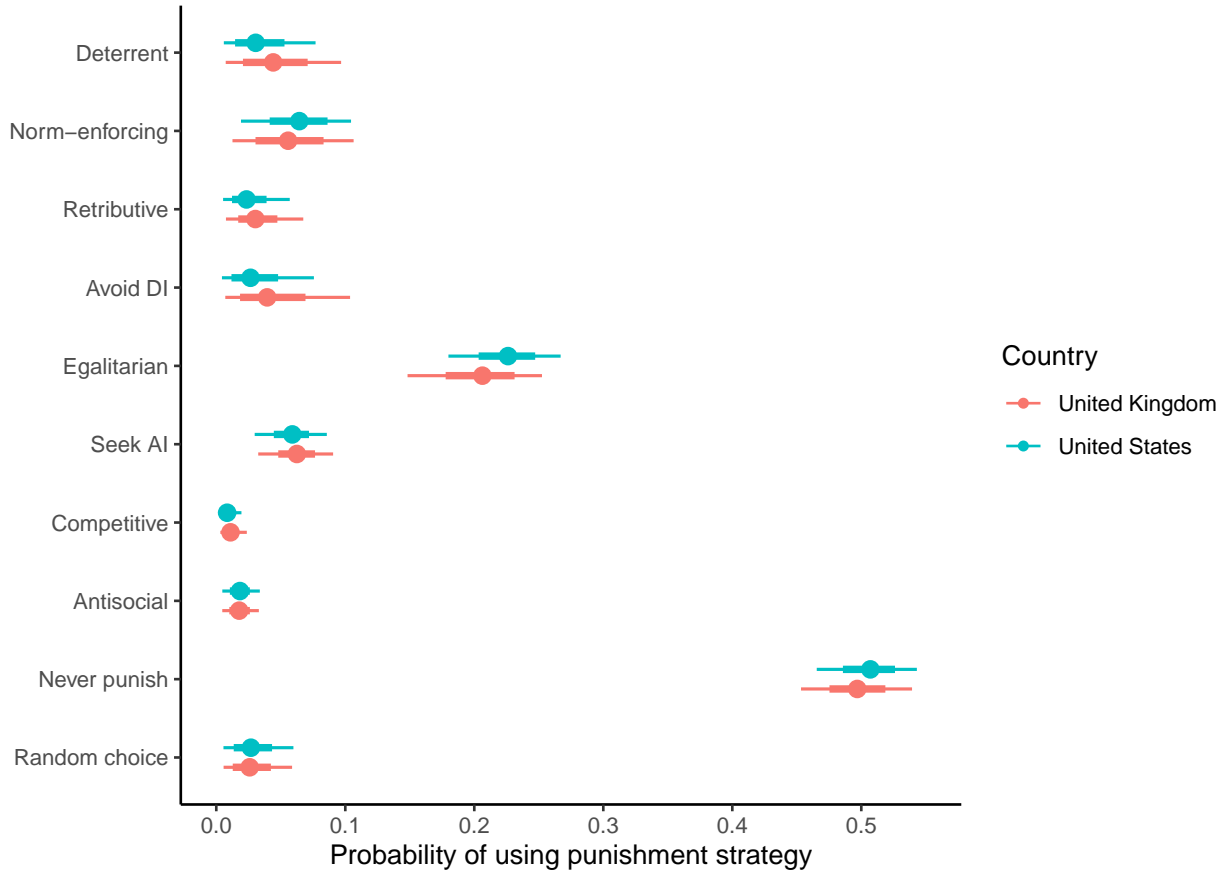


Figure 3. Posterior estimates of the probabilities of following different punishment strategies from the Bayesian latent state model. The model assumes an implementation error rate of 5%. Points represent posterior medians, line ranges represent 50% and 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.

Next, we explored which traits predicted adherence to different punishment strategies. To answer this question, we included variables capturing demographics, personality, social preferences, political views, and religious views as predictors in our Bayesian latent state model. We included each variable in a separate model, predicting all ten punishment strategies (the nine from Table 1, plus the ‘random choice’ strategy) simultaneously.

Demographic variables tended to be unrelated to strategy usage: age and gender did

not predict adherence to a particular punishment strategy (Supplementary Figures S2 and S3). In the United States, the never punish strategy was slightly more common among participants lower in socio-economic status (median posterior slope = -0.20, 95% CI [-0.38 -0.03]) but this effect was small.

Conversely, personality and social preferences were linked to variation in punishment strategies. When including the Big-6 personality dimensions and Social Value Orientation (SVO) in the model, we found associations with the egalitarian, never punish, and random choice strategies, with small-to-medium effect sizes (Figure 4). More prosocial participants (those with higher SVO scores) were more likely to follow the egalitarian and the never punish strategies, while those with lower SVO scores were more likely to enact the random choice strategy. The personality dimensions of honesty-humility and openness to experience were both positively associated with following the never punish strategy, while extraversion negatively predicted this strategy. The effects were mostly similar across countries, but occasionally differed: for example, in the United States, but not in the United Kingdom, honesty-humility was positively associated with following the egalitarian strategy and negatively associated with following the random choice strategy. Overall, the same pattern of results emerged when analysing the full dataset without exclusions (Supplementary Figure S4).

Political and religious variables were also associated with punishment strategy (Figure 5). These effects were small-to-medium in size and tended to be more pronounced in the United States. Controlling for Social Dominance Orientation, American participants higher in Right Wing Authoritarianism were more likely to follow the strategies avoiding disadvantageous inequity and seeking advantageous inequity. Participants who stated that they would like to “bring those below them [on the socio-economic status ladder] up a peg” were more likely to follow the egalitarian strategy, while American participants higher in Social Dominance Orientation, Right Wing Authoritarianism, and believing that God controls events in the world were less likely to follow the egalitarian strategy. In general,

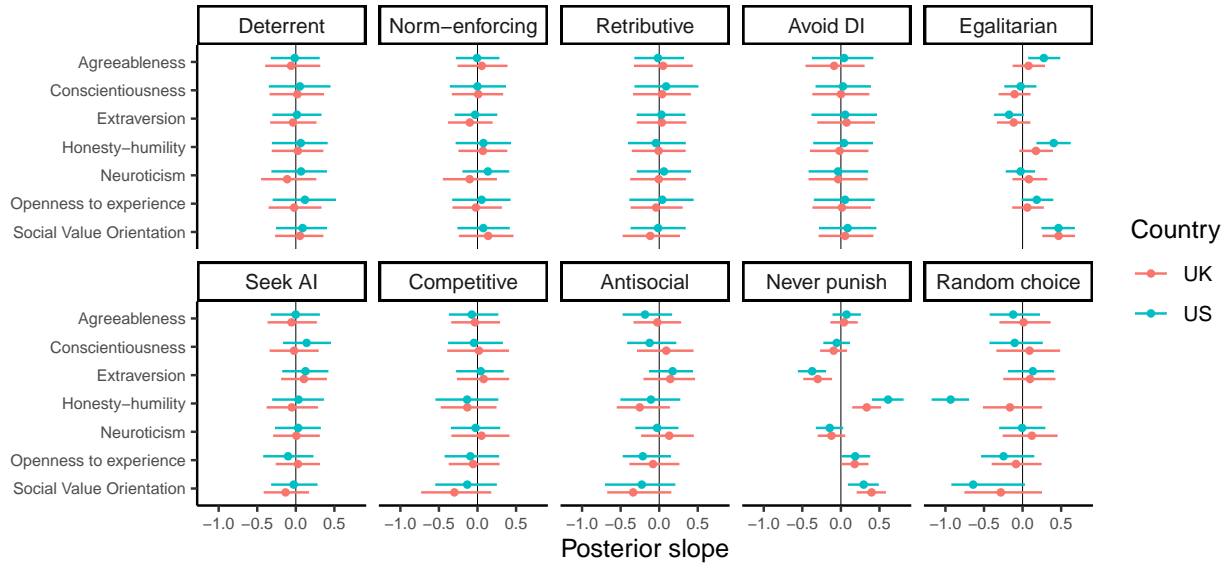


Figure 4. Posterior slopes from Bayesian latent state models including Big-6 personality dimensions and Social Value Orientation. Each row represents a separate model. All models assume an implementation error rate of 5%. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.

religious and conservative participants were less likely to follow the never punish strategy. This general pattern of results was replicated with the full dataset (Supplementary Figure S5).

Finally, we asked whether participants had insight into their own punishment strategy. In other words, could participants self-report the strategy that they were following during the games? To answer this question, we included participants' responses to post-game questions about their strategy as predictors in the model. As before, each predictor was included in a separate model, predicting all ten strategies simultaneously.

In general, we found that self-reported strategy usage was positively associated with the behavioural strategy that participants employed, with effects ranging from small to large in size (see Supplementary Figures S6 and S7 for the distribution of responses to self-report questions). Figure 6 shows the relationships between self-report questions and the different punishment strategies, highlighting the combinations where the question matched the behavioural strategy. We found positive relationships between the self-report

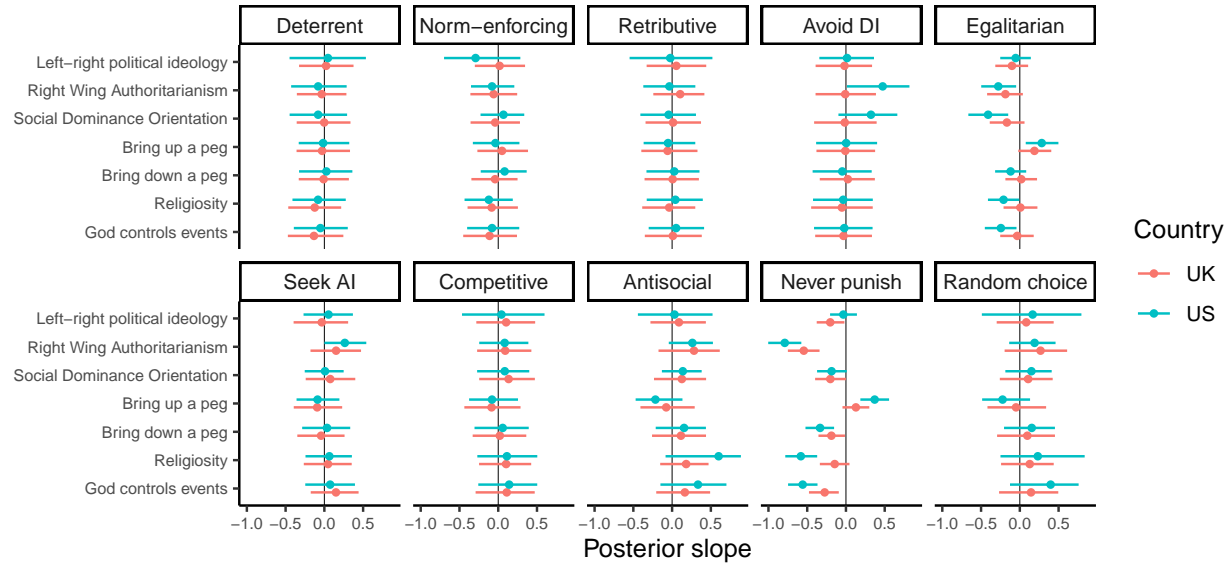


Figure 5. Posterior slopes from Bayesian latent state models including political ideology, views about social inequality, and religiosity. Each row represents a separate model aside from Social Dominance Orientation and Right Wing Authoritarianism, which control for one another within the same model. All models assume an implementation error rate of 5%. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.

questions and strategy usage for the norm-enforcing, egalitarian, seek advantageous inequity, never punish, and random choice strategies. The 95% credible intervals for other estimates included zero, though these estimates often trended in a positive direction. The same pattern of results was found when analysing the full dataset without exclusions (Supplementary Figure S8).

Discussion

Using a suite of economic games measuring punishment in different situations, we have shown that punishment does not serve just one function, but instead is a flexible tool that can be and is used for different purposes (6). While some use punishment to enforce norms of cooperation, others use it to reduce or even create inequality between individuals. We found that people's punishment strategy can, to some extent, be predicted by individual differences in personality, social preferences, and political and religious views. Moreover,

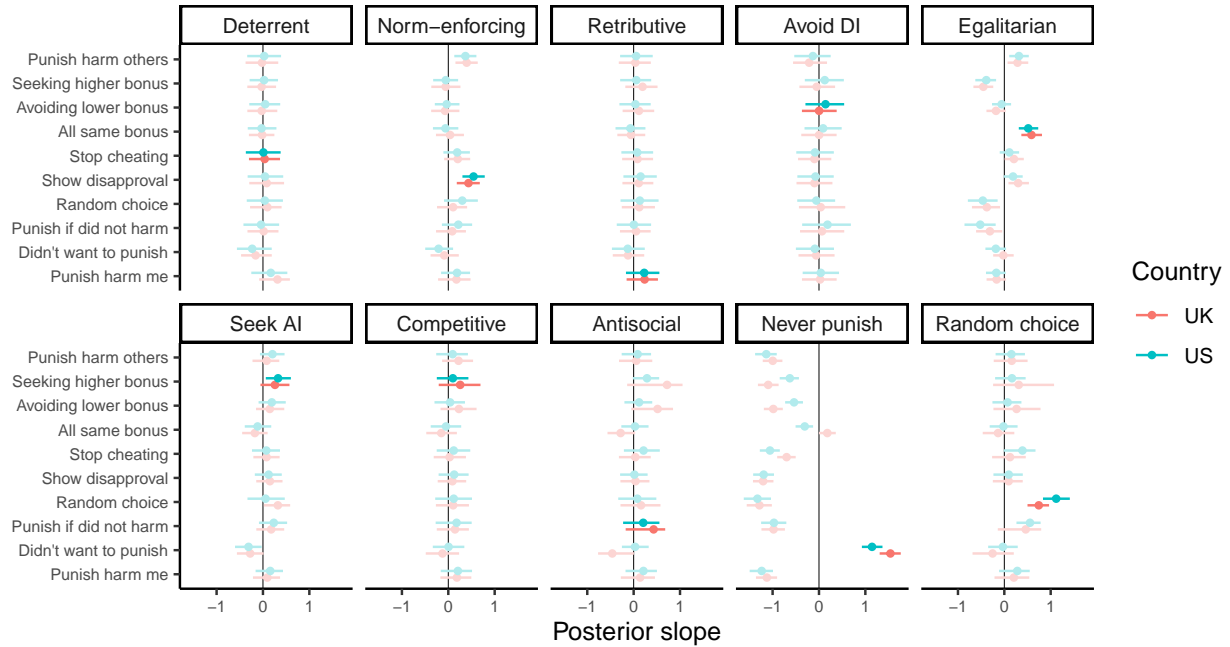


Figure 6. Posterior slopes from models including self-reported strategy usage. Each row represents a separate model. All models assume an implementation error rate of 5%. Highlighted estimates represent combinations where the self-report question matched the behavioural strategy. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.

contrary to the view that people are often unable to articulate the reasons for their punitive behaviour (28), people seem to have some degree of insight into the strategy they are using. Despite small differences, these general patterns replicated in samples from both the United Kingdom and the United States, providing further confidence in the results.

Among the punitive strategies, the most common were particularly sensitive to inequality in payoffs, either from a self-referential perspective (i.e., avoid disadvantageous inequity) or more generally (i.e., egalitarian). This is in line with previous studies which have highlighted inequity aversion as an important motivation for punishment in economic games (29–31, 37, 48, 50). Personality, social preferences, and political and religious views were related to these strategies. Traits associated with other-regarding concern, such as SVO and honesty-humility, predicted following the egalitarian strategy. By contrast, religious and conservative individuals, including individuals higher in SDO and RWA, were

less likely to follow the egalitarian strategy and more likely to follow the avoid disadvantageous inequity strategy, especially in the United States. Authoritarian participants in the United States were also more likely to actively seek advantageous inequity with their punishment. These findings provide insight into the motives that may have driven higher levels of peer punishment among conservatives in previous work (41).

Behaviour shaping strategies, such as deterrence and norm-enforcement, were less common than strategies sensitive to inequality in our set of games. This was reflected both in participants' elicited punishment behaviour (Figure 3) and in their self-reports of their own strategy (Supplementary Figures S6 and S7). Although our design did not explicitly allow for behaviour shaping as the interactions were all one-shot, we did manipulate whether the target's stealing behaviour was intentional or not (Game C), an approach which has been used in previous vignette studies to identify behaviour shaping motives (28). The lower prevalence of behaviour shaping strategies in our study is consistent with prior work showing that punishment often continues to be used when behaviour shaping is impossible, such as when the target will never find out that they have been punished (34) or on the last round of repeated interactions (33). We found that demographic, personality, political, and religious variables tended to be unrelated to behaviour shaping strategies, though this could reflect low power to detect such associations given the low prevalence of these strategies in our sample. We also found that participants accurately reported using the norm-enforcing strategy, but not the deterrent strategy. This finding is in line with previous research showing that people struggle to accurately report the deterrent motivations for their punitive behaviour (28).

We defined antisocial punishment as occurring when a participant exclusively punished those who did not steal. There were no participants in our sample who followed this strategy, although some players did punish non-stealing co-players. Harming non-stealing individuals was also consistent with the competitive strategies, which did appear in our sample albeit at low frequencies. The fact that no participants in our sample

247 punished non-stealing across all games suggests that the antisocial punishment that has
248 been observed in other studies (17, 51) is not aimed at harming cooperators specifically, as
249 has been previously suggested (17). Instead, antisocial punishment is more likely to be
250 motivated by improving one's relative position, which is in line with work showing that
251 antisocial punishment disappears when relative payoffs cannot be changed (e.g., when
252 punishment is only available with a 1:1 fee-fine ratio) (35).

253 The fact that people use punishment for many different reasons poses problems for
254 the way that punishment is operationalised in classic behavioural economic game studies.
255 In these studies, a common assumption is that participants punish to change the behaviour
256 of cheats (11, 12). But in reality, people may be choosing the punishment option to achieve
257 a variety of different goals. This has implications for how people respond to being targeted
258 by punishers in these games. Targets of punishment in these studies may know that
259 punishment reflects different motives and can respond accordingly. For example, if targets
260 interpret punishment as serving a competitive motive, it may elicit retaliation rather than
261 encourage cooperation (6, 10, 52). As punishers' motives must be inferred (and such
262 inferences likely depend on character traits of the target, as well as the context in which
263 punishment occurs), there is likely to be some variation and error in attributing motives to
264 punishers. To the extent that inferred motives affect target responses, this might help to
265 explain the mixed findings in the field as to whether punishment actually motivates
266 cheating targets to cooperate in the future (6).

267 It is striking that the most common strategy in our dataset was to never punish. This
268 is partly because punishment in these games imposes an economic cost for no tangible
269 benefit. If the fee-fine ratio had been lower, such that it was cheaper to punish, we may
270 have seen more punishment from participants. Indeed, 72% of participants following the
271 never punish strategy positively stated that they didn't want to pay to reduce anyone's
272 bonus but would have done so if it were free. But the frequency of the never punish
273 strategy perhaps also reflects a more general aversion to peer punishment, an aversion that

has been highlighted in both WEIRD (Western, educated, industrialised, rich, and democratic) samples (53, 54) and in small-scale societies (55). One reason that people may be averse to peer punishment is that, because it is a fundamentally harmful act, punishment can reflect badly on the punisher and people might therefore refrain from punishing others to avoid reputational damage (4). People frequently avoid taking actions that could harm their reputation, even when they don't know if reputation is at stake (as in the one-shot anonymous settings used here). Another reason that people may be averse to peer punishment is that it can trigger retaliation (6). This may be especially likely in situations that lack clear institutional norms to legitimise punishment, such as our economic games. As with reputation damage, people might abstain from peer punishment to avoid retaliation, regardless of whether retaliation is actually possible. By contrast, institutionalised punishment in small-scale societies often functions to compensate victims while limiting the potential for feuds and cycles of retaliation (56, 57). Future research should uncover whether people are more willing to punish in these conventionalised contexts (e.g., see ref (58)).

One potential limitation with our design is that some strategies required more punishment than others, meaning that some strategies were more “expensive” to implement. For example, the competitive strategy required punishment in ten of twelve decisions, compared to the avoid disadvantageous inequity strategy which required only one instance of punishment (Table 1). We employed the strategy method to deal with this, calculating participant payoffs from a randomly chosen game instead of accumulating the costs across all games. This mitigates the concern that the total cost is the most salient feature determining participants' decisions. But even accounting for this, one could still argue that the difference in overall costs explains why, for example, the competitive strategy is less common in our dataset than the avoid disadvantageous inequity strategy. We do not think this feature of our design is a concern, however, for a number of reasons. First, we are interested in measuring strategies underlying *costly* punishment. Some of

these strategies, by their very nature, will be more costly to implement than others. This is reflected in our design. It would not make theoretical sense to use an alternative design where more or less punitive strategies are manipulated to cost participants the same amount. Second, when we asked participants whether they would have punished if it were free to do so, we found low agreement with this statement (Supplementary Figure S6), suggesting that participants were not particularly sensitive to the costs of punishment. In line with this, when we plot the frequencies of different strategies against their expected costs, we find that cost is not a perfect predictor of strategy frequency: many “cheap” strategies are rare in our data and some “expensive” strategies are quite common (e.g., always punish; Supplementary Figure S9). Third, this argument implies that the high frequency of the never punish strategy is merely an artefact of our design, since it is the only strategy that does not cost anything to implement. But as we have discussed, many other studies have found similar aversions to costly punishment in the lab (29, 31, 33, 44, 48) and in the real world (53–55) suggesting that this result is not an artefact of our design.

In sum, we have shown that while many people choose not to punish peers, those who do are motivated by a variety of different concerns, including behaviour shaping, egalitarianism, and competition. Much like the observed variation in human social learning strategies (39), humans thus also exhibit variation in their punishment strategies. These individual differences map onto personality dimensions, social preferences, political and religious views, and self-reports of behaviour. We hope that future work will continue to unpack the multifaceted nature of human punishment.

Methods

Ethical approval

This research has been approved by the UCL Department of Psychology Ethics Committee (project: 3720/002) and ratified by the University of Auckland Human

Participants Ethics Committee. The study was performed in accordance with all the relevant guidelines and regulations. Informed consent was obtained from all participants prior to the study and participation was voluntary.

Pre-registration

We pre-registered the study on the Open Science Framework before collecting data in the United Kingdom (11th November 2022; <https://osf.io/k75fc>). We submitted another pre-registration before collecting data in the United States (20th June 2023; <https://osf.io/q4hdy>). In the pre-registrations, we outlined our study design, exclusion criteria, and analysis plan. As the study was exploratory, we did not pre-register any explicit hypotheses. We did not deviate from the pre-registrations.

Exclusion criteria

We pre-registered that we would exclude participants who failed any of the attention checks, sped through the surveys (i.e., two standard deviations below the median duration), or flatlined (i.e., provided identical responses to matrix questions). We also stated that we would exclude data for particular games if participants failed the comprehension question for that game. We followed our pre-registered plan of conducting analyses with and without these exclusions (analyses without exclusions are reported in the Supplementary Material).

Participants

We collected a representative sample of 1019 participants from the United Kingdom through the online platform Prolific (<https://www.prolific.com/>). All of these participants completed the economic games and 973 returned to complete the follow-up survey a week later (95% retention rate). After exclusions, we were left with 1014 participants overall (see Supplementary Figure S10 for sample characteristics).

We later collected a representative sample of 1005 participants from the United States through Prolific. All of these participants completed the economic games and 957 returned to complete the follow-up survey (95% retention rate). After exclusions, we were left with 996 participants overall (see Supplementary Figure S11 for sample characteristics).

Materials

Economic games. In the first part of the study, participants completed six economic games, each with slight variations. In all games, there are multiple players and the participant takes the role of P1. P2 either (a) steals £0.20 from another player and adds it to their payoff or (b) does nothing. For each of these cases, participants are asked whether they would like to pay money to reduce P2’s payoff. Games A-E have two players and Game F has three players.

The six games are as follows (variations bolded; see Figure 1 for a visual representation of the games):

1. *Game A (Advantageous Inequity)*. P1 starts with £0.70 and P2 starts with £0.10. P2 is given the option to either steal £0.20 from P1 or do nothing. P1 can then pay £0.10 to reduce P2’s payoff by £0.30.
2. *Game B (Equal)*. P1 starts with £0.70 and P2 starts with **£0.30**. P2 is given the option to either steal £0.20 from P1 or do nothing. P1 can then pay £0.10 to reduce P2’s payoff by £0.30.
3. *Game C (Computer)*. P1 starts with £0.70 and P2 starts with £0.30. Participants are told that “**the computer will decide**” whether P2 steals £0.20 from P1 or does nothing. P1 can then pay £0.10 to reduce P2’s payoff by £0.30.
4. *Game D (1:1 Fee-Fine)*. P1 starts with £0.70 and P2 starts with £0.30. P2 is given the option to either steal £0.20 from P1 or do nothing. P1 can then pay £0.10 to reduce P2’s payoff by **£0.10**.

5. *Game E (Disadvantageous Inequity)*. P1 starts with £0.70 and P2 starts with **£0.50**.

P2 is given the option to either steal £0.20 from P1 or do nothing. P1 can then pay £0.10 to reduce P2's payoff by £0.30.

6. *Game F (Third-Party)*. P1 starts with £1.00, P2 and P3 start with £0.70. P2 is

given the option to either steal £0.20 **from P3** or do nothing. P1 can then pay £0.10 to reduce P2's payoff by £0.30.

Our punishment strategies make different behavioural predictions across these different games (Table 1). The deterrent strategy punishes whenever it is intentionally personally harmed (when the target steals in Games A, B, D, and E). The norm-enforcing strategy is the same, except that it also punishes when *others* are intentionally harmed (when the target steals in Game F). The retributive strategy punishes whenever it is personally harmed, whether the harm was intentional or not (when the target steals in Games A-E). The avoid disadvantageous inequity strategy only punishes when it currently experiences disadvantageous inequity and can alleviate it by punishing (when the target steals in Game E). The egalitarian strategy punishes when punishment makes payoffs for all players more equal than they currently are (when the target steals in Games E and F). The seek advantageous inequity strategy punishes when it currently does not experience advantageous inequity and punishment would produce advantageous inequity for itself (when the target steals in Games B and C). The competitive strategy punishes whenever the fee-fine ratio is 1:3, as it can improve its relative position by punishing, but does not punish when the fee-fine ratio is 1:1 (Game D). The antisocial strategy punishes whenever the target does not steal across all games. The never punish strategy never punishes in any game.

For each game, participants saw the game instructions and answered a comprehension question before providing their decisions. After completing all the games, participants were asked to give an open-ended response explaining their behaviour in the games, and then

responded to several slider questions capturing the different reasons for their decisions (for full wordings, see Supplementary Table S2).

Survey questions. In a follow-up survey, we collected the following data on participants (for wordings of all questions, see Supplementary Table S3):

- *Demographics.* In the survey, we collected information on participants' education level and self-reported socio-economic status (MacArthur ladder (59)). We also collected additional demographic data from Prolific (e.g., age, gender, student status).
- *Personality.* We used the Mini-IPIP scale (60) to measure the Big 6 personality dimensions of agreeableness ($\alpha = 0.83$), conscientiousness ($\alpha = 0.75$), extraversion ($\alpha = 0.83$), honesty-humility ($\alpha = 0.77$), openness to experience ($\alpha = 0.79$), and neuroticism ($\alpha = 0.79$). Four items were used for each personality dimension.
- *Social Value Orientation.* We used the Social Value Orientation Slider Measure to measure other-regarding preferences (61). Across fifteen items, participants made decisions on how to allocate different amounts of money between themselves and another anonymous individual. From these decisions, we calculated participants' Social Value Orientation "angle" as a measure of their other-regarding preference, following the steps outlined in ref (61).
- *Political ideology.* We included several measures of political ideology, including left-right conservatism, Social Dominance Orientation (62) ($\alpha = 0.91$; eight items), and Right Wing Authoritarianism (63) ($\alpha = 0.82$; six items). We also probed participants' views on social inequality by asking them whether they would like to bring people above (below) them on the MacArthur socio-economic status ladder down (up) a peg or two.
- *Religious views.* We asked participants how religious they consider themselves and whether they believe that God or another spiritual non-human entity controls the events in the world (64).

Procedure

We began data collection in the United Kingdom on 28th November 2022, with participants returning to complete the follow-up survey on 5th December 2022. We then ran a second wave of data collection in the United States on 20th June 2023, with participants returning to complete the follow-up survey on 27th June 2023. Our surveys were designed through the online survey platform Qualtrics (<https://www.qualtrics.com/>).

In the initial games survey, participants completed all six economic games in a random order, with punishment decisions (whether to punish a stealing target and whether to punish a target who did nothing) randomised within games. Responses to comprehension questions suggested that participants understood the six economic games (Supplementary Table S4). In order to partially mitigate the fact that some punishment strategies predict more punishment than others and are thus more “expensive” to implement (Table 1), we used the strategy method to incentivise the economic games, choosing a random game to determine bonus payments rather than summing participants’ earnings across all games. After all games, 62% of participants stated that they believed that their decisions had real consequences for others.

In the follow-up survey, participants completed blocks of questions on demographics, personality, Social Value Orientation, political ideology, and religious views in a random order, with questions randomised within blocks. A random decision from the Social Value Orientation Slider Measure was chosen to determine bonus payment.

Participants were paid £1.80 for completing the games survey, plus a bonus payment from the six economic games (between £0.40 – £0.70 depending on their decisions). Participants were paid £1.50 for completing the follow-up survey, plus a bonus payment from the Social Value Orientation Slider Measure (between £0.50 – £0.85 depending on their decisions).

Statistical analysis

We pre-registered that we would use a Bayesian latent state model to infer unobserved punishment strategies from the observed data [for a similar version of this model, see ref (65)]. In this model, participants i in countries c make binary punishment decisions across twelve decisions j . We assume that the probability of the observed data $y_{i,j}$ is the weighted average of the probability of the observed data conditional on each of the ten punishment strategies s . From this logic, the model estimates the probability of each strategy p_s . The full model is as follows:

$$\begin{aligned}
 y_{i,j} &\sim \text{Bernoulli}(\theta_j) \\
 \theta_j &= \sum_{s=1}^{10} p_s \text{Pr}(\text{punish}|s, j) \\
 p &= \text{softmax}(\alpha_{c[i]}) \\
 \alpha_{s,c} &\sim \text{Normal}(0, 1)
 \end{aligned} \tag{1}$$

The conditional probabilities $\text{Pr}(\text{punish}|s, j)$ are hard coded in the model as outlined in Table 1. We incorporate an implementation error rate δ into these conditional probabilities by coding green ticks in Table 1 with a conditional probability of $1 - \delta$ and coding red crosses with a conditional probability of $0 + \delta$. We set δ to 0.05 in all models, which is similar to its value when we estimate it as a free parameter in an additional model (median posterior $\delta = 0.03$, 95% CI [0.00 0.06]; Supplementary Figure S12). The random choice strategy is consistently coded with a conditional probability of $\frac{1}{2}$ across all decisions.

To include a categorical predictor in the model, we estimate a different $\alpha_{s,c}$ for each categorical level. To include a continuous predictor x in the model, we include a slope β in the linear model for p :

$$\begin{aligned}
y_{i,j} &\sim \text{Bernoulli}(\theta_j) \\
\theta_j &= \sum_{s=1}^{10} p_s \text{Pr}(\text{punish}|s, j) \\
p &= \text{softmax}(\alpha_{c[i]} + \beta_{c[i]} x_i) \\
\alpha_{s,c} &\sim \text{Normal}(0, 1) \\
\beta_{s,c} &\sim \text{Normal}(0, 0.2)
\end{aligned} \tag{2}$$

These models control for multiple comparisons across strategies by estimating the effects of the predictor on all strategies simultaneously.

We estimated the posterior distributions of these models using Hamiltonian Monte Carlo as implemented in Stan version 2.26.1 (66). We ran each model for 2000 samples, with 1000 warmup samples. R-hat values and effective sample sizes suggested that all models converged normally. Trace plots are reported in Supplementary Figure S13.

We validated the model by simulating observed data ($n = 100$) from a known frequency of strategies. The model was successfully able to recover the known frequency of strategies from the simulated data (Supplementary Figure S14).

Reproducibility

All data and code are accessible on GitHub:
<https://github.com/ScottClaessens/punishStrategies>. All analyses were conducted in R version 4.2.1 (67). Visualisations were created with the *ggplot2* (68) and *cowplot* (69) R packages. We used the *targets* (70) R package to create a reproducible data analysis pipeline and the *papaja* (71) R package to reproducibly generate the manuscript.

Acknowledgements

This work was supported by a Royal Society of New Zealand Catalyst Leaders Grant to Q.D.A and N.R. (ref: ILFUOA2002).

Author Contributions

All authors (S.C., Q.D.A., N.R.) conceptualised the research, designed the study, and developed the surveys. N.R. conducted data collection on Prolific. S.C. conducted all analyses and visualisation of the data. All authors wrote the manuscript.

Competing Interests

The authors declare no competing interests.

Data Availability

All data used in this study are publicly available on GitHub:
<https://github.com/ScottClaessens/punishStrategies>

Code Availability

All code to reproduce the analyses in this study are publicly available on GitHub:
<https://github.com/ScottClaessens/punishStrategies>

References

1. T. H. Clutton-Brock, G. A. Parker, Punishment in animal societies. *Nature* **373**, 209–216 (1995).
2. M. dos Santos, D. J. Rankin, C. Wedekind, The evolution of punishment through reputation. *Proceedings of the Royal Society B: Biological Sciences* **278**, 371–377 (2011).
3. M. dos Santos, D. J. Rankin, C. Wedekind, Human cooperation based on punishment reputation. *Evolution* **67**, 2446–2450 (2013).
4. N. J. Raihani, R. Bshary, The reputation of punishers. *Trends in Ecology & Evolution* **30**, 98–103 (2015).
5. E. Ostrom, *Governing the commons: The evolution of institutions for collective action* (Cambridge University Press, 1990).
6. N. J. Raihani, R. Bshary, Punishment: One tool, many uses. *Evolutionary Human Sciences* **1**, e12 (2019).
7. D. Balliet, L. B. Mulder, P. A. M. Van Lange, Reward, punishment, and cooperation: A meta-analysis. *Psychological Bulletin* **137**, 594–615 (2011).
8. D. Balliet, P. A. M. V. Lange, Trust, punishment, and cooperation across 18 societies: A meta-analysis. *Perspectives on Psychological Science* **8**, 363–379 (2013).
9. A. Chaudhuri, Sustaining cooperation in laboratory public goods experiments: A selective survey of the literature. *Experimental Economics* **14**, 47–83 (2011).
10. A. Dreber, D. G. Rand, D. Fudenberg, M. A. Nowak, Winners don't punish. *Nature* **452**, 348–351 (2008).
11. E. Fehr, S. Gächter, Cooperation and punishment in public goods experiments. *American Economic Review* **90**, 980–994 (2000).
12. E. Fehr, S. Gächter, Altruistic punishment in humans. *Nature* **415**, 137–140 (2002).

13. J. Henrich *et al.*, Costly punishment across human societies. *Science* **312**, 1767–1770 (2006).
14. N. Nikiforakis, H.-T. Normann, A comparative statics analysis of punishment in public-good experiments. *Experimental Economics* **11**, 358–369 (2008).
15. N. J. Raihani, A. Thornton, R. Bshary, Punishment and cooperation in nature. *Trends in Ecology & Evolution* **27**, 288–295 (2012).
16. D. J.-F. de Quervain *et al.*, The neural basis of altruistic punishment. *Science* **305**, 1254–1258 (2004).
17. B. Herrmann, C. Thöni, S. Gächter, Antisocial punishment across societies. *Science* **319**, 1362–1367 (2008).
18. C. F. Camerer, *Behavioral game theory: Experiments in strategic interaction* (Russell Sage Foundation, 2003).
19. S. Bowles, H. Gintis, *A cooperative species: Human reciprocity and its evolution* (Princeton University Press, 2013).
20. R. Boyd, H. Gintis, S. Bowles, P. J. Richerson, The evolution of altruistic punishment. *Proceedings of the National Academy of Sciences* **100**, 3531–3535 (2003).
21. M. Chudek, J. Henrich, Culture-gene coevolution, norm-psychology and the emergence of human prosociality. *Trends in Cognitive Sciences* **15**, 218–226 (2011).
22. J. Henrich, *The secret of our success: How culture is driving human evolution, domesticating our species, and making us smarter* (Princeton University Press, 2017).
23. A. W. Delton, M. M. Krasnow, The psychology of deterrence explains why group membership matters for third-party punishment. *Evolution and Human Behavior* **38**, 734–743 (2017).
24. E. Fehr, I. Schurtenberger, Normative foundations of human cooperation. *Nature Human Behaviour* **2**, 458–468 (2018).

25. S. Mathew, R. Boyd, Punishment sustains large-scale cooperation in prestate warfare. *Proceedings of the National Academy of Sciences* **108**, 11375–11380 (2011).
26. S. Mathew, R. Boyd, The cost of cowardice: Punitive sentiments towards free riders in Turkana raids. *Evolution and Human Behavior* **35**, 58–64 (2014).
27. P. Richerson *et al.*, Cultural group selection plays an essential role in explaining human cooperation: A sketch of the evidence. *Behavioral and Brain Sciences* **39**, e30 (2016).
28. K. M. Carlsmith, J. M. Darley, P. H. Robinson, Why do we punish? Deterrence and just deserts as motives for punishment. *Journal of Personality and Social Psychology* **83**, 284–299 (2002).
29. N. J. Raihani, K. McAuliffe, Human punishment is motivated by inequity aversion, not a desire for reciprocity. *Biology Letters* **8**, 802–804 (2012).
30. C. T. Dawes, J. H. Fowler, T. Johnson, R. McElreath, O. Smirnov, Egalitarian motives in humans. *Nature* **446**, 794–796 (2007).
31. J. E. Bone, N. J. Raihani, Human punishment is motivated by both a desire for revenge and a desire for equality. *Evolution and Human Behavior* **36**, 323–330 (2015).
32. J. M. Walker, M. A. Halloran, Rewards and sanctions and the provision of public goods in one-shot settings. *Experimental Economics* **7**, 235–247 (2004).
33. P. Barclay, N. Raihani, Partner choice versus punishment in human prisoner’s dilemmas. *Evolution and Human Behavior* **37**, 263–271 (2016).
34. M. J. Crockett, Y. Özdemir, E. Fehr, The value of vengeance and the demand for deterrence. *Journal of Experimental Psychology: General* **143**, 2279–2286 (2014).
35. K. Sylwester, B. Herrmann, J. J. Bryson, Homo homini lupus? Explaining antisocial punishment. *Journal of Neuroscience, Psychology, and Economics* **6**, 167–188 (2013).
36. E. Fehr, U. Fischbacher, Third-party punishment and social norms. *Evolution and Human Behavior* **25**, 63–87 (2004).

37. T. Johnson, C. T. Dawes, J. H. Fowler, R. McElreath, O. Smirnov, The role of egalitarian motives in altruistic punishment. *Economics Letters* **102**, 192–194 (2009).
38. J. H. Fowler, T. Johnson, O. Smirnov, Egalitarian motive and altruistic punishment. *Nature* **433**, E1 (2005).
39. L. Molleman, P. Van den Berg, F. J. Weissing, Consistent individual differences in human social learning strategies. *Nature Communications* **5**, 3570 (2014).
40. I. Thielmann, G. Spadaro, D. Balliet, Personality and prosocial behavior: A theoretical framework and meta-analysis. *Psychological Bulletin* **146**, 30–90 (2020).
41. W. Hofmann, M. J. Brandt, D. C. Wisneski, B. Rockenbach, L. J. Skitka, Moral punishment in everyday life. *Personality and Social Psychology Bulletin* **44**, 1697–1711 (2018).
42. R. E. Nisbett, T. D. Wilson, Telling more than we can know: Verbal reports on mental processes. *Psychological Review* **84**, 231–259 (1977).
43. P. Barclay, Reputational benefits for altruistic punishment. *Evolution and Human Behavior* **27**, 325–344 (2006).
44. T. Batistoni, P. Barclay, N. J. Raihani, Third-party punishers do not compete to be chosen as partners in an experimental game. *Proceedings of the Royal Society B* **289**, 20211773 (2022).
45. J. J. Jordan, M. Hoffman, P. Bloom, D. G. Rand, Third-party punishment as a costly signal of trustworthiness. *Nature* **530**, 473–476 (2016).
46. J. J. Jordan, D. G. Rand, Third-party punishment as a costly signal of high continuation probabilities in repeated games. *Journal of Theoretical Biology* **421**, 189–202 (2017).
47. J. J. Jordan, D. G. Rand, Signaling when no one is watching: A reputation heuristics account of outrage and punishment in one-shot anonymous interactions. *Journal of*

Personality and Social Psychology **118**, 57–88 (2020).

48. P. Deutchman, M. Bračić, N. J. Raihani, K. McAuliffe, Punishment is strongly motivated by revenge and weakly motivated by inequity aversion. *Evolution and Human Behavior* **42**, 12–20 (2021).

49. J. Marczyk, Human punishment is not primarily motivated by inequality. *PLoS ONE* **12**, e0171298 (2017).

50. J. E. Bone, K. McAuliffe, N. J. Raihani, Exploring the motivations for punishment: Framing and country-level effects. *PLoS One* **11**, e0159769 (2016).

51. A. Pleasant, P. Barclay, Why hate the good guy? Antisocial punishment of high cooperators is greater when people compete to be chosen. *Psychological Science* **29**, 868–876 (2018).

52. N. Nikiforakis, Punishment and counter-punishment in public good games: Can we really govern ourselves? *Journal of Public Economics* **92**, 91–112 (2008).

53. L. Balafoutas, N. Nikiforakis, B. Rockenbach, Direct and indirect punishment among strangers in the field. *Proceedings of the National Academy of Sciences* **111**, 15924–15927 (2014).

54. L. Balafoutas, N. Nikiforakis, B. Rockenbach, Altruistic punishment does not increase with the severity of norm violations in the field. *Nature Communications* **7**, 13327 (2016).

55. N. Baumard, Has punishment played a role in the evolution of cooperation? A critical review. *Mind & Society* **9**, 171–192 (2010).

56. L. Fitouchi, M. Singh, Punitive justice serves to restore reciprocal cooperation in three small-scale societies. *Evolution and Human Behavior* **44**, 502–514 (2023).

57. M. Singh, Z. H. Garfield, Evidence for third-party mediation but not punishment in Mentawai justice. *Nature Human Behaviour* **6**, 930–940 (2022).

58. L. Molleman, F. Kölle, C. Starmer, S. Gächter, People prefer coordinated punishment in cooperative interactions. *Nature Human Behaviour* **3**, 1145–1153 (2019).
59. N. E. Adler, E. S. Epel, G. Castellazzo, J. R. Ickovics, Relationship of subjective and objective social status with psychological and physiological functioning: Preliminary data in healthy, White women. *Health Psychology* **19**, 586–592 (2000).
60. C. G. Sibley *et al.*, The Mini-IPIP6: Validation and extension of a short measure of the Big-Six factors of personality in New Zealand. *New Zealand Journal of Psychology (Online)* **40**, 142 (2011).
61. R. O. Murphy, K. A. Ackermann, M. J. J. Handgraaf, Measuring social value orientation. *Judgment and Decision Making* **6**, 771–781 (2011).
62. A. K. Hoet *al.*, The nature of social dominance orientation: Theorizing and measuring preferences for intergroup inequality using the new sdo7 scale. *Journal of Personality and Social Psychology* **109**, 1003 (2015).
63. B. Bizumic, J. Duckitt, Investigating right wing authoritarianism with a Very Short Authoritarianism Scale. *Journal of Social and Political Psychology* **6**, 129–150 (2018).
64. K. Laurin, A. F. Shariff, J. Henrich, A. C. Kay, Outsourcing punishment to God: Beliefs in divine control reduce earthly punishment. *Proceedings of the Royal Society B: Biological Sciences* **279**, 3272–3281 (2012).
65. R. McElreath, *Statistical rethinking: A Bayesian course with examples in R and Stan, 2nd edition*, 2nd Ed. (CRC Press, 2020).
66. Stan Development Team, RStan: The R interface to Stan (2020).
67. R Core Team, *R: A language and environment for statistical computing* (R Foundation for Statistical Computing, 2022).
68. H. Wickham, *ggplot2: Elegant graphics for data analysis* (Springer-Verlag New York, 2016).

69. C. O. Wilke, *cowplot: Streamlined plot theme and plot annotations for 'ggplot2'* (2020).
70. W. M. Landau, The targets R package: A dynamic Make-like function-oriented pipeline toolkit for reproducibility and high-performance computing. *Journal of Open Source Software* **6**, 2959 (2021).
71. F. Aust, M. Barth, *papaja: Prepare reproducible APA journal articles with R Markdown* (2022).

Supplementary Material

Why do people punish? Evidence for a range of strategic concerns

Scott Claessens¹, Quentin D. Atkinson¹, & Nichola Raihani^{1,2}

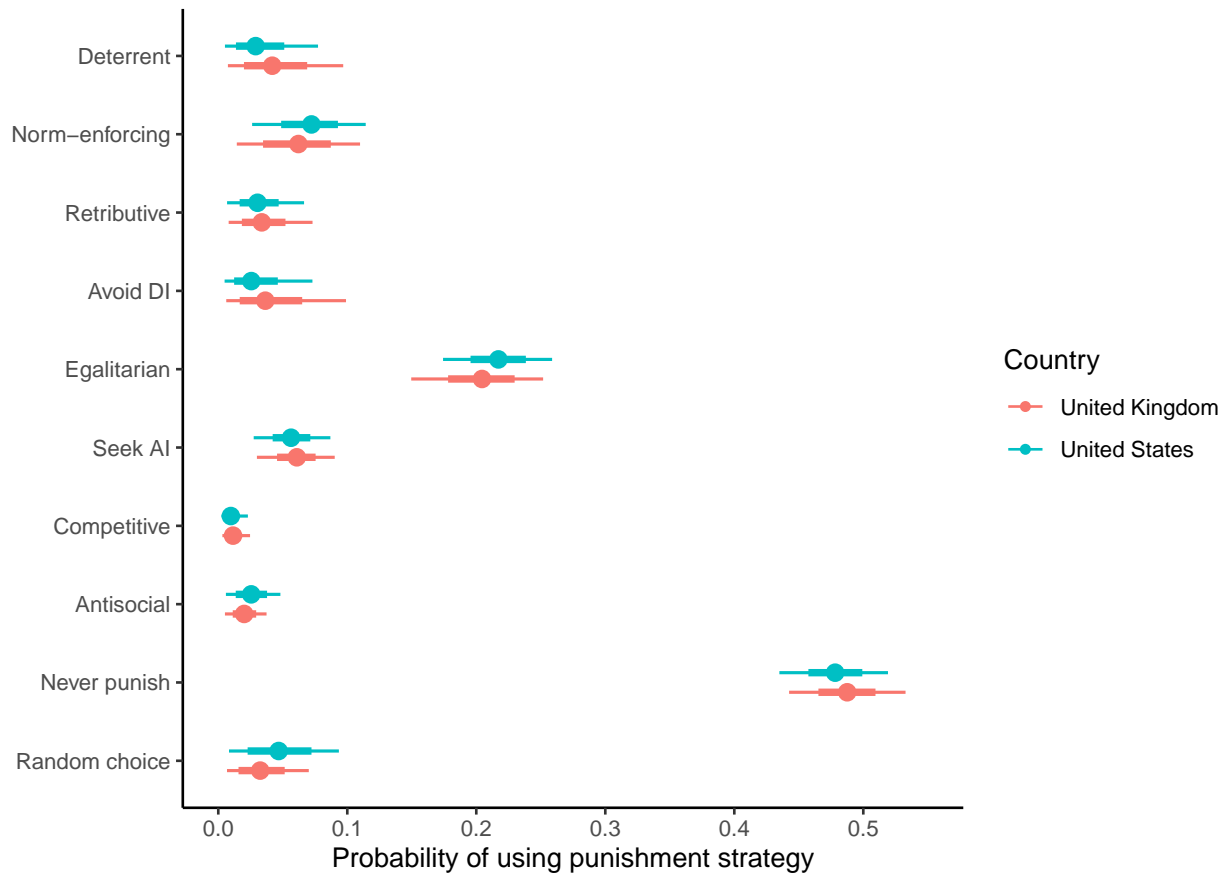
¹ School of Psychology, University of Auckland, Auckland, New Zealand

² Department of Experimental Psychology, University College London, London, United Kingdom

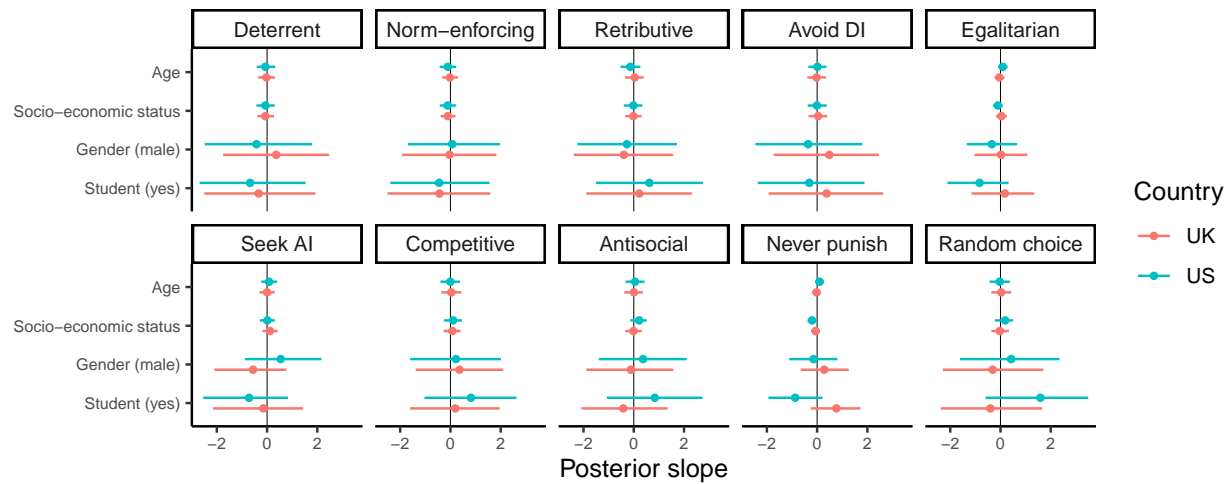
Contents

Supplementary Figures S1-S14	2
Supplementary Tables S1-S4	16

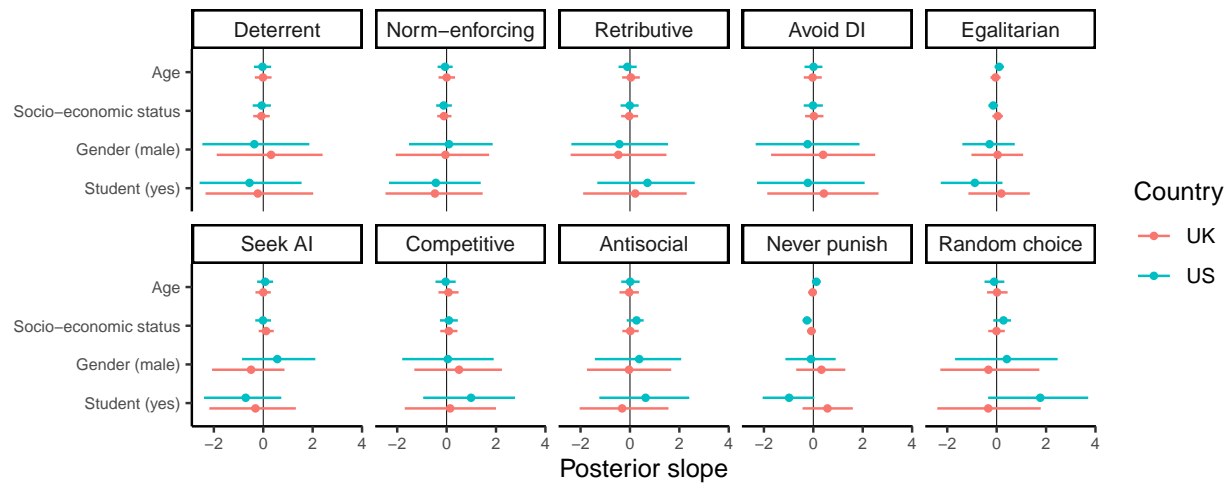
Supplementary Figures



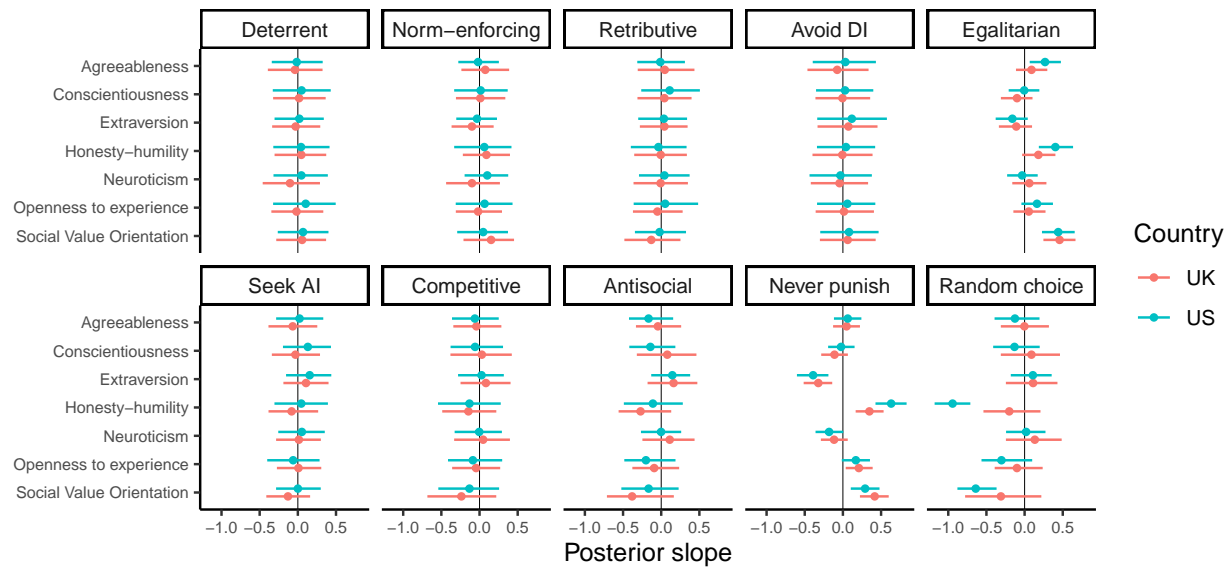
Supplementary Figure S1. Posterior estimates of the probabilities of following different punishment strategies from the Bayesian latent state model fitted to the full dataset without pre-registered exclusions. The model assumes an implementation error rate of 5%. Figure 3 in the main text shows the same result, but from a model fitted to the reduced dataset with pre-registered exclusions. Points represent posterior medians, line ranges represent 50% and 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



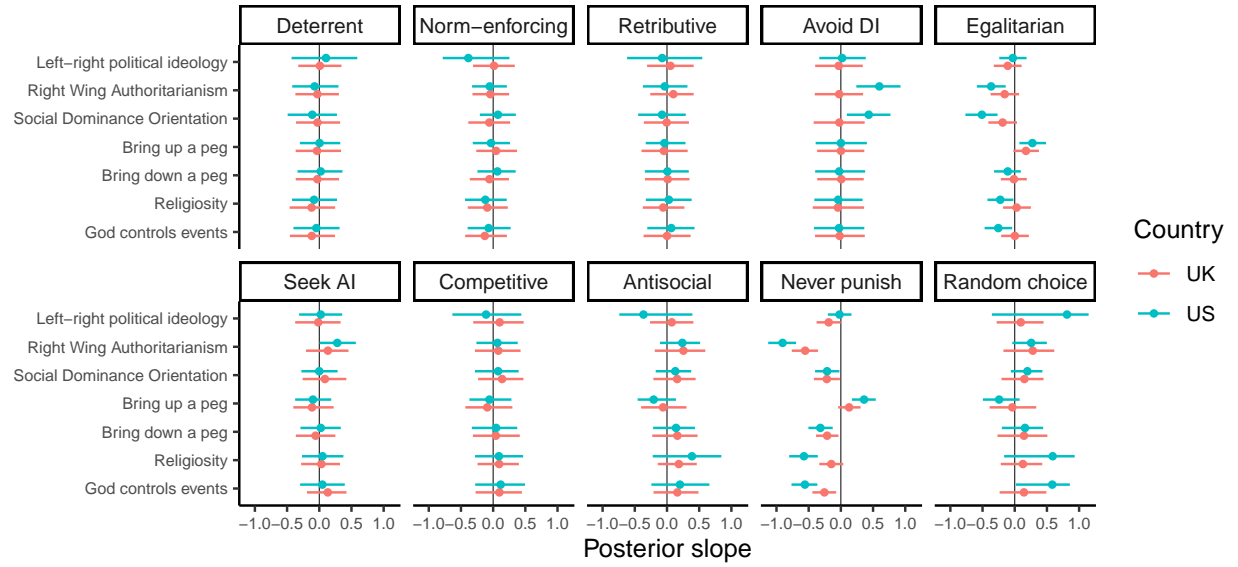
Supplementary Figure S2. Posterior slopes from models including age, socio-economic status, gender, and student status, fitted to the subsetting dataset with pre-registered exclusions. Each row represents a separate model. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



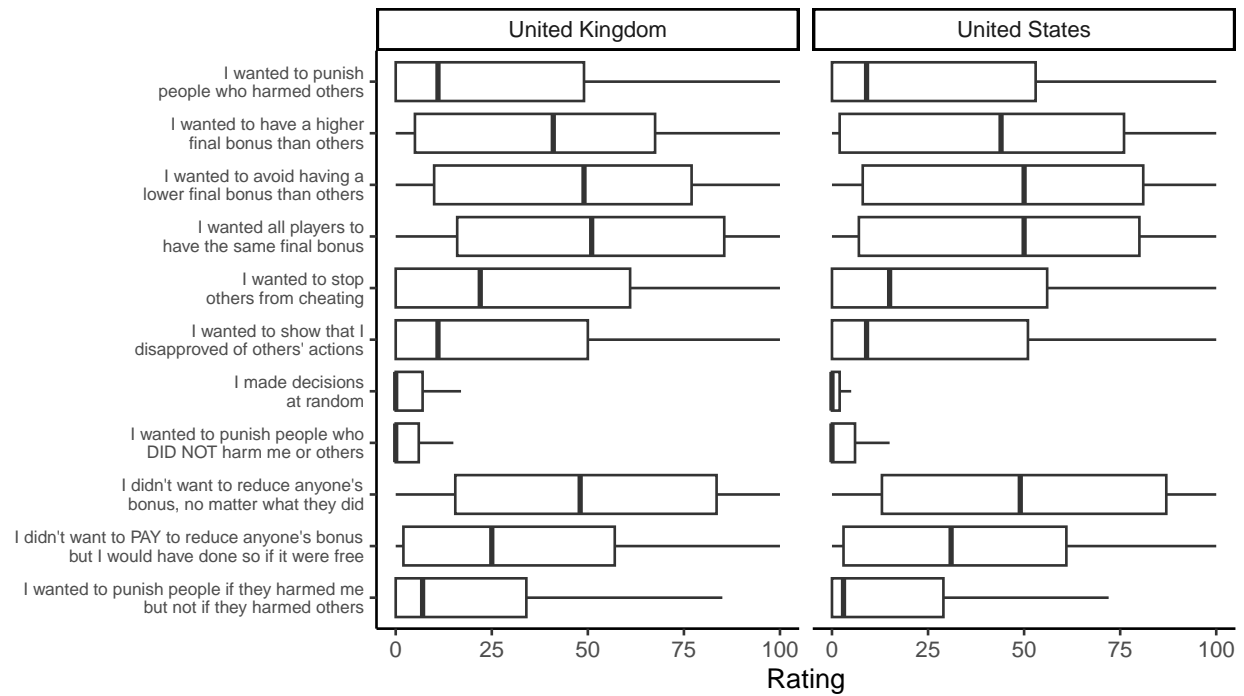
Supplementary Figure S3. Posterior slopes from models including age, socio-economic status, gender, and student status, fitted to the full dataset without pre-registered exclusions. Each row represents a separate model. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



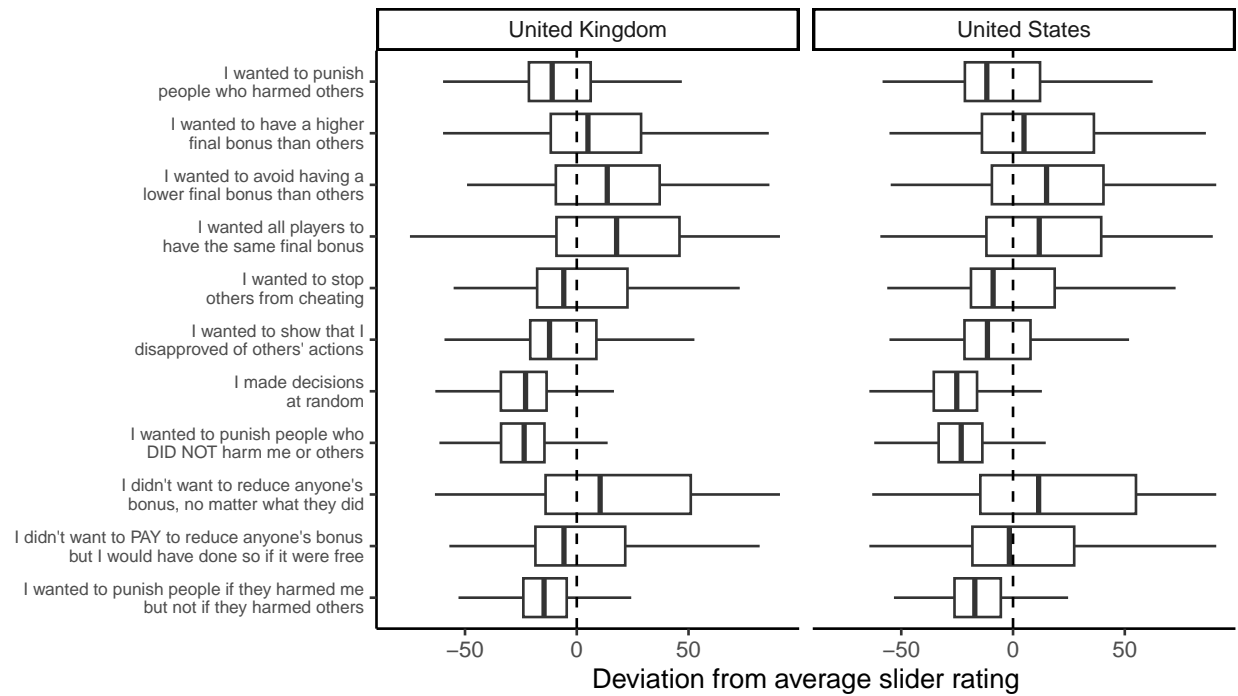
Supplementary Figure S4. Posterior slopes from models including Big-6 personality dimensions and Social Value Orientation, fitted to the full dataset without pre-registered exclusions. Each row represents a separate model. Figure 4 in the main text shows the same results, but from models fitted to the reduced dataset with pre-registered exclusions. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



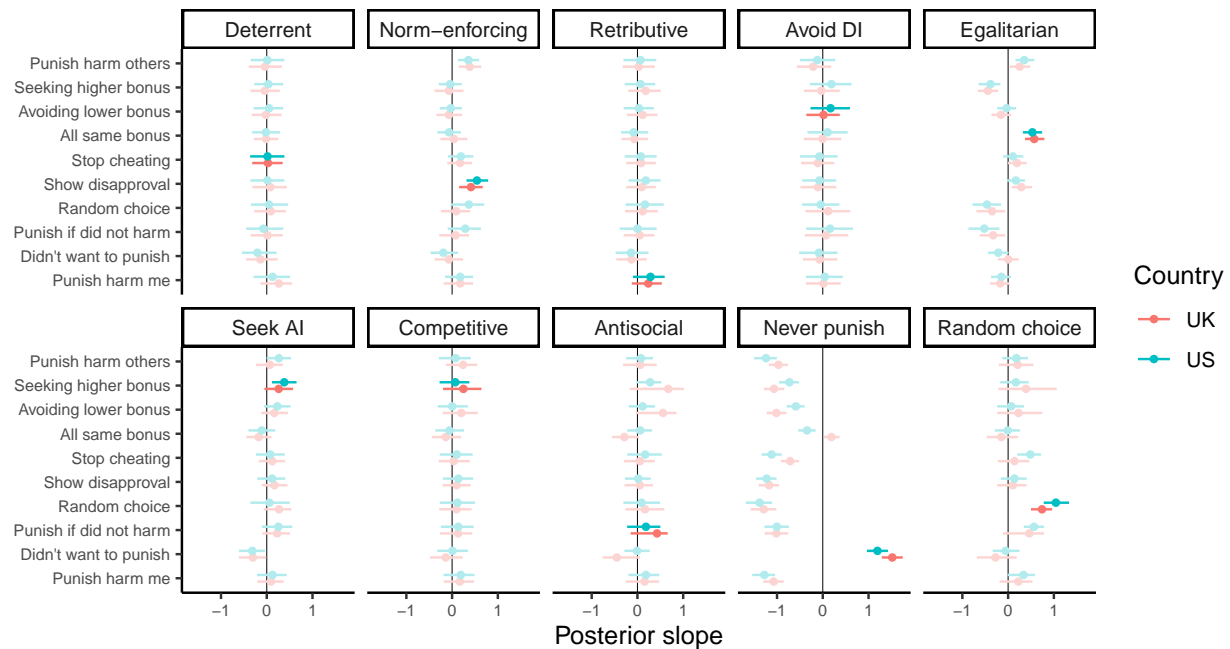
Supplementary Figure S5. Posterior slopes from models including political ideology, views about social inequality, and religiosity, fitted to the full dataset without pre-registered exclusions. Each row represents a separate model aside from Social Dominance Orientation and Right Wing Authoritarianism, which control for one another within the same model. Figure 5 in the main text shows the same results, but from models fitted to the reduced dataset with pre-registered exclusions. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



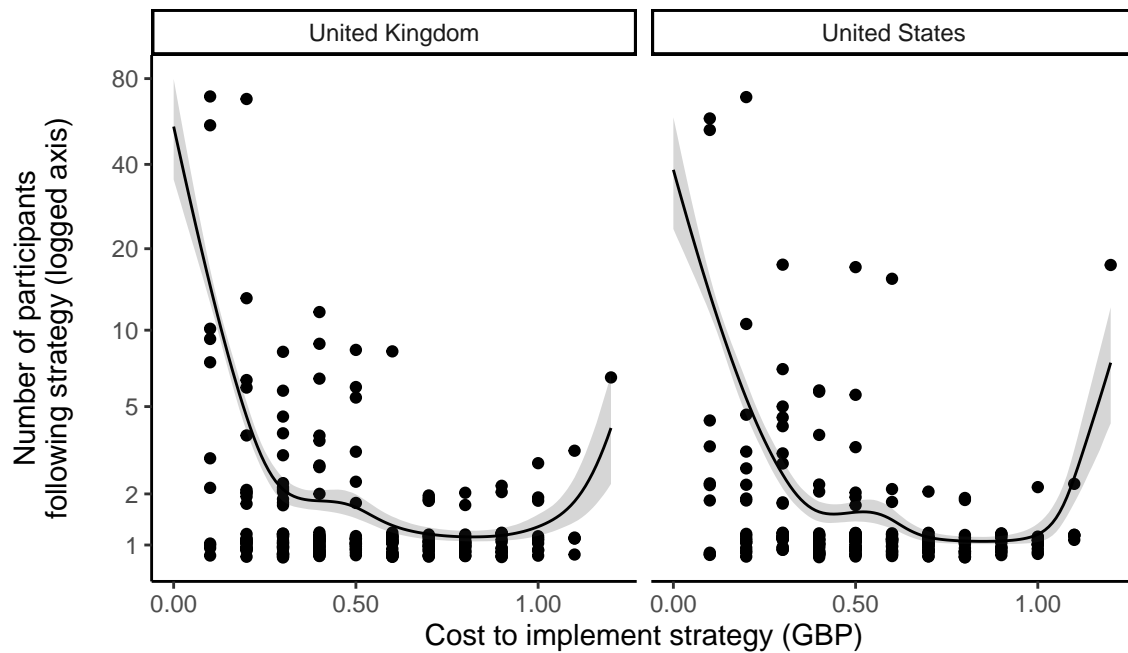
Supplementary Figure S6. Boxplots showing the distribution of responses to each self-report question about the reasons for participants' behaviour in the games. Boxplots represent medians and interquartile ranges.



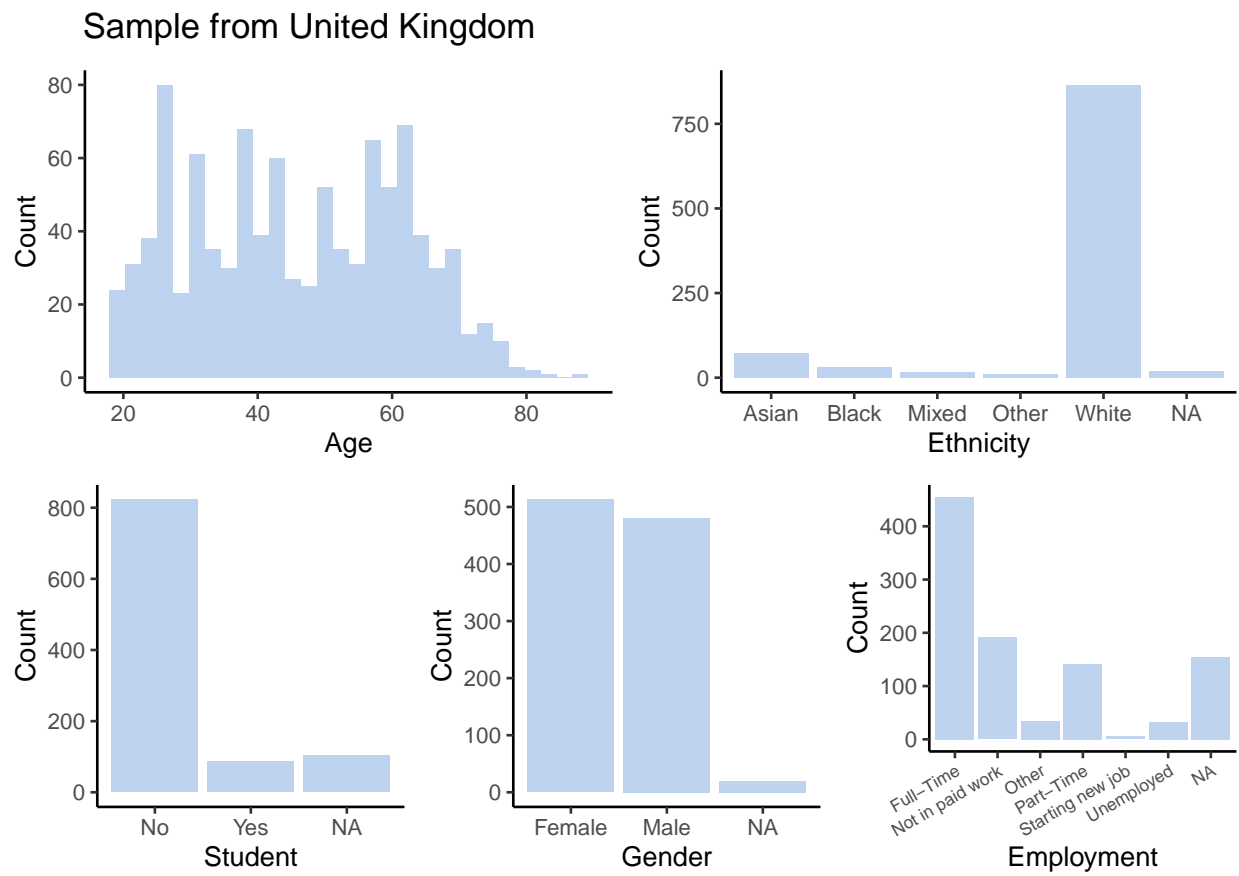
Supplementary Figure S7. Boxplots showing the distribution of responses to each self-report question about the reasons for participants' behaviour in the games, presented as deviations from participants' average rating across all questions. Boxplots represent medians and interquartile ranges.



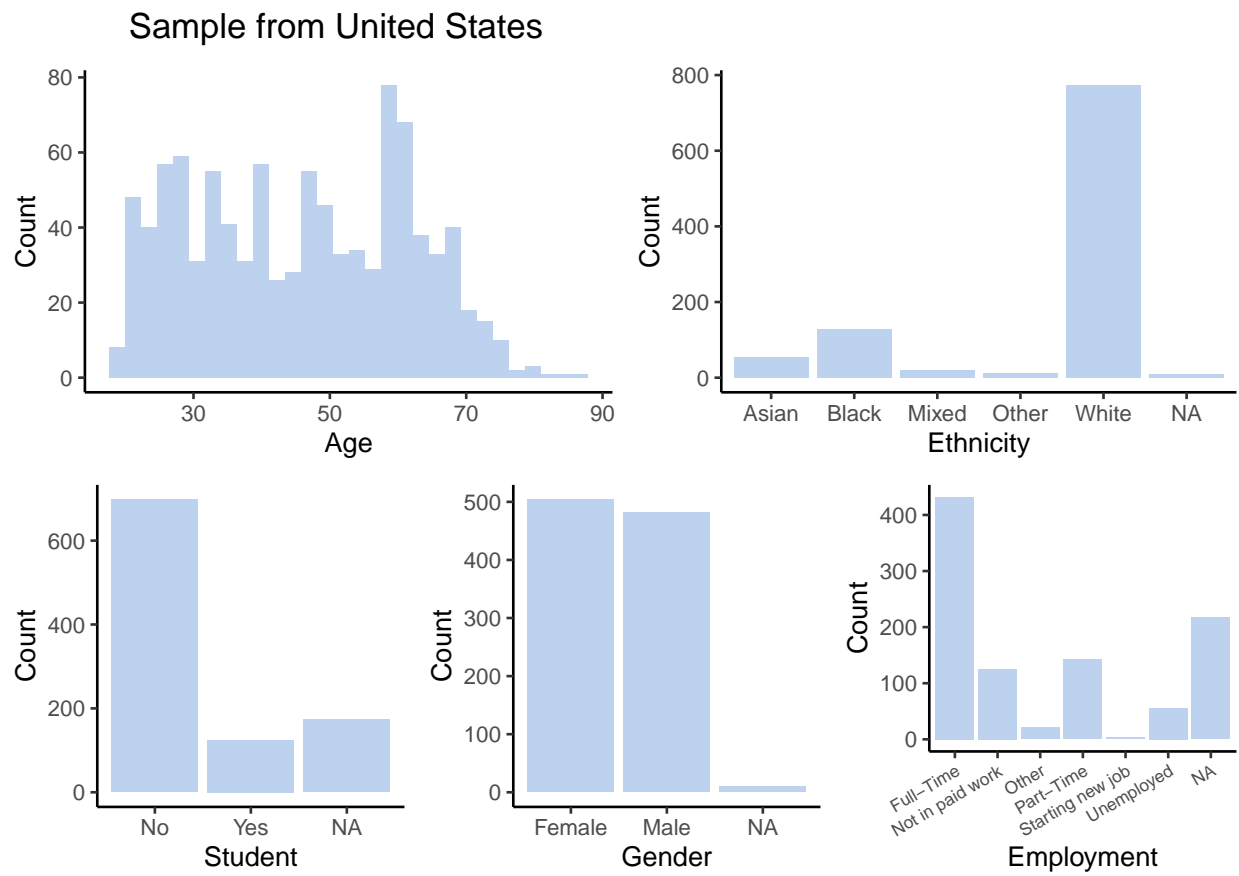
Supplementary Figure S8. Posterior slopes from models including self-reported strategy usage, fitted to the full dataset without pre-registered exclusions. Each row represents a separate model. Highlighted estimates represent combinations where the self-report slider matched the behavioural strategy. Figure 6 in the main text shows the same results, but from models fitted to the reduced dataset with pre-registered exclusions. Points represent posterior medians, line ranges represent 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



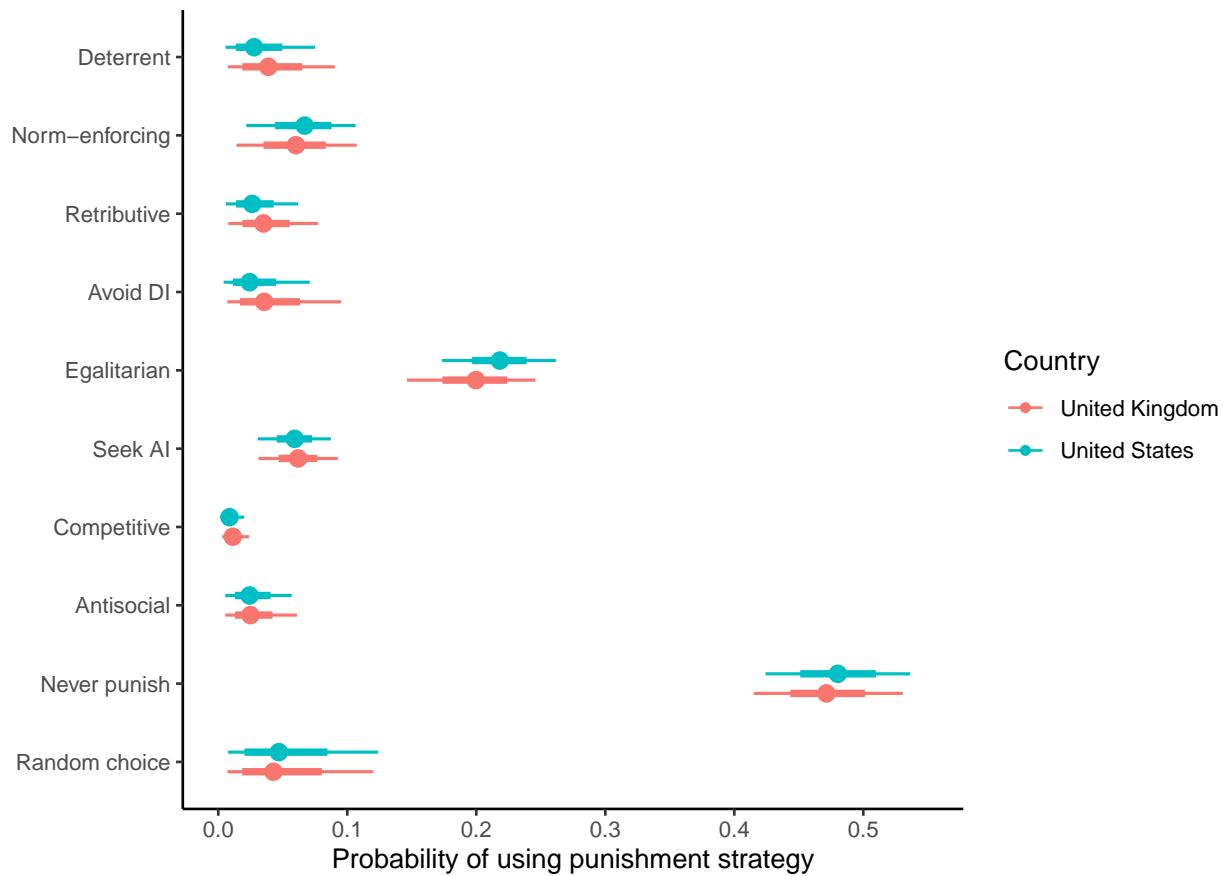
Supplementary Figure S9. The relationships between strategy frequencies and the overall costs of strategies across all twelve punishment decisions, in both countries. Each point is a unique strategy that appears in our dataset at least once (for ease of presentation, the “never punish” strategy is excluded). Lines and shaded areas represent posterior predictions from splines fitted to the trend in each country separately.



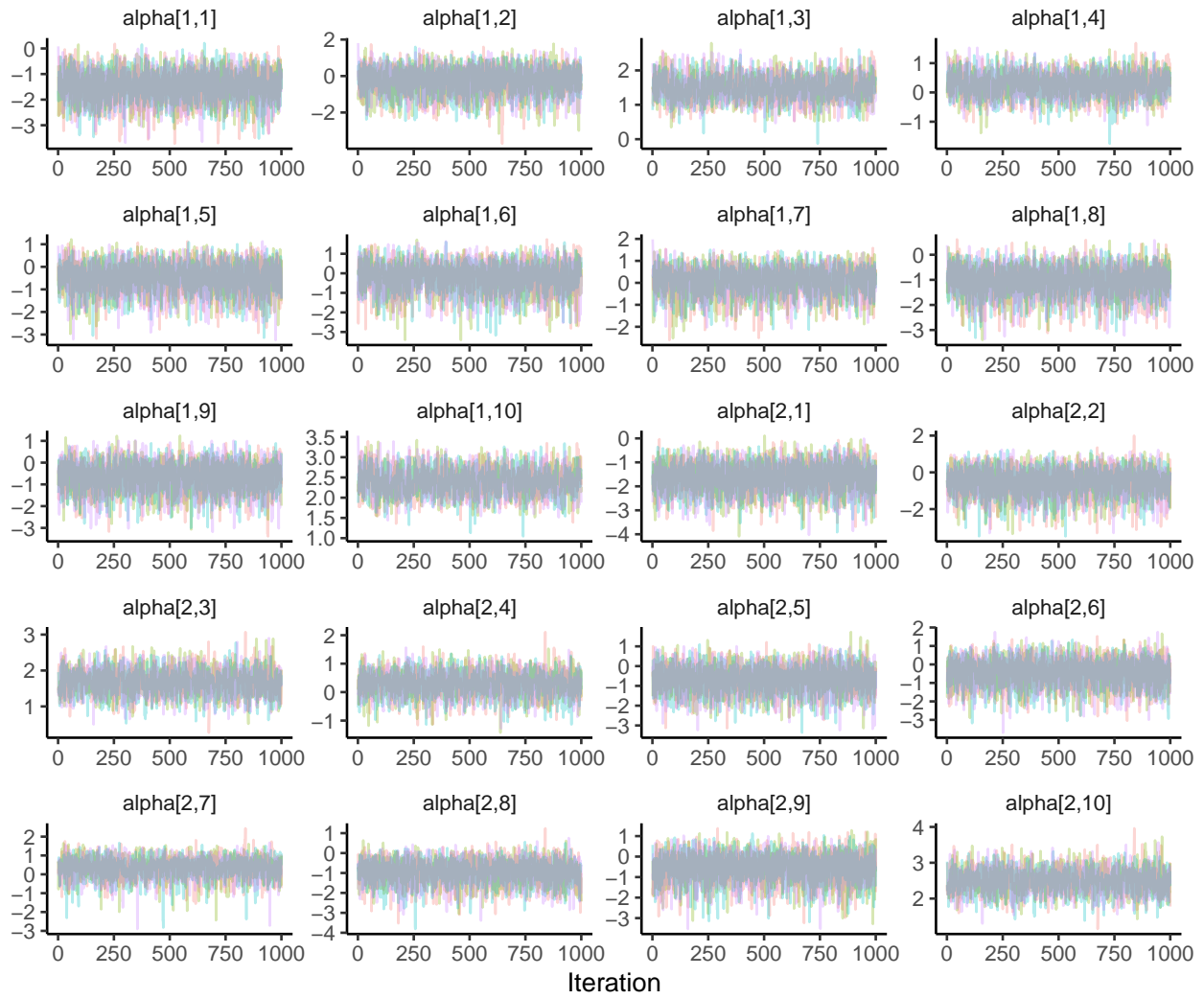
Supplementary Figure S10. Sample characteristics in the United Kingdom.



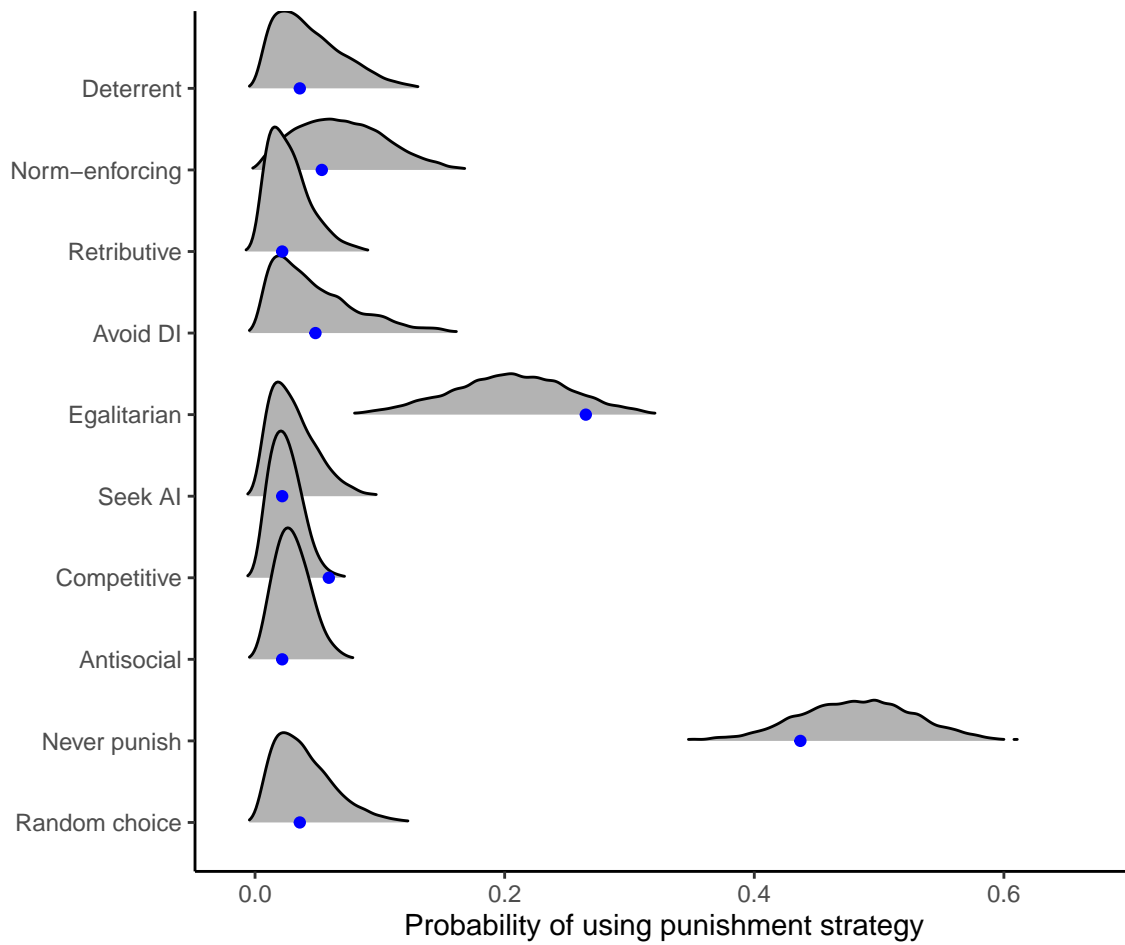
Supplementary Figure S11. Sample characteristics in the United States.



Supplementary Figure S12. Posterior estimates of the probabilities of following different punishment strategies from the Bayesian latent state model that estimates the implementation error rate as a free parameter. The model estimated the implementation error rate to be 0.03 (95% CI [0.00 0.06]). Points represent posterior medians, line ranges represent 50% and 95% credible intervals. AI = advantageous inequity, DI = disadvantageous inequity.



Supplementary Figure S13. MCMC trace plots for parameter values from the Bayesian latent state model fitted to data with exclusions. These trace plots suggest that the different chains mixed well and the model converged normally.



Supplementary Figure S14. Results of Bayesian latent state model fitted to simulated data ($n = 100$) with known strategy frequencies in the population. Blue points represent known strategy frequencies, grey densities represent posterior estimates of strategy frequencies. AI = advantageous inequity, DI = disadvantageous inequity.

Supplementary Tables

Supplementary Table S1

Counts and proportions of the 25 most common patterns of punitive behaviour across all twelve decisions, split by country. Binary strings represent punishment (1) or no punishment (0) in each decision, aligning with the order of game decision columns in Table 1.

Pattern	Explanation	United Kingdom (N = 1014)		United States (N = 996)	
		N	Prop	N	Prop
000000000000	<i>Never punish strategy (exact)</i>	426	0.420	447	0.449
000000001000	<i>Avoid DI strategy (exact)</i>	67	0.066	62	0.062
000000001010	<i>Egalitarian strategy (exact)</i>	65	0.064	71	0.071
000000000010	Punish when steal in Game F	55	0.054	49	0.049
001000001000	Punish when steal in Games B and E	14	0.014	11	0.011
101000001010	Punish when steal in Games A, B, E, and F	11	0.011	4	0.004
100000000000	Punish when steal in Game A	10	0.010	2	0.002
000000100000	Punish when steal in Game D	9	0.009	3	0.003
001000001010	Punish when steal in Games B, E, and F	9	0.009	17	0.017
101000101000	<i>Deterrent strategy (exact)</i>	9	0.009	6	0.006
101010101010	Punish when steal in all games	9	0.009	15	0.015
101000101010	<i>Norm-enforcing strategy (exact)</i>	8	0.008	16	0.016
001000000000	Punish when steal in Game B	7	0.007	4	0.004
001010101000	Punish when steal in Games B, C, D, and E	7	0.007	0	0.000
100000001000	Punish when steal in Games A and E	6	0.006	5	0.005
101000001000	Punish when steal in Games A, B, and E	6	0.006	7	0.007
101010101000	<i>Retributive strategy (exact)</i>	6	0.006	5	0.005
111111111111	Always punish	6	0.006	16	0.016
000000101000	Punish when steal in Games D and E	5	0.005	2	0.002
000000101010	Punish when steal in Games D, E, and F	5	0.005	3	0.003
101010001010	Punish when steal in all games except Game D	5	0.005	2	0.002
001000101000	Punish when steal in Games B, D, and E	4	0.004	2	0.002
001000101010	Punish when steal in Games B, D, E, and F	4	0.004	6	0.006
101000000000	Punish when steal in Games A and B	4	0.004	2	0.002
101010001000	Punish when steal in Games A, B, C, and E	4	0.004	0	0.000

Supplementary Table S2

Wordings for 11 self-report slider questions asking participants to report the reasons for their behaviour in the six games. Participants were prompted with the following text: “We would now like you to answer a few questions about your main motivation in the games. Please answer truthfully - there is no right or wrong answer and your first answer is probably best. Please rate the extent to which the following statements apply to your decisions to reduce or not to reduce other players’ bonuses in the games.”

Slider	Wording
1	I wanted to punish people who harmed others
2	I wanted to have a higher final bonus than others
3	I wanted to avoid having a lower final bonus than others
4	I wanted all players to have the same final bonus
5	I wanted to stop others from cheating
6	I wanted to show that I disapproved of others’ actions
7	I made decisions at random
8	I wanted to punish people who DID NOT harm me or others
9	I didn’t want to reduce anyone’s bonus, no matter what they did
10	I didn’t want to PAY to reduce anyone’s bonus but I would have done so if it were free
11	I wanted to punish people if they harmed me but not if they harmed others

Wordings for survey questions in the study.

[illegible]

Table S3 continued

Measure	Wording	Scale
Right Wing Authoritarianism	It's great that many young people today are prepared to defy authority (reversed)	1-9
	What our country needs most is discipline, with everyone following our leaders in unity	1-9
	God's laws about abortion, pornography, and marriage must be strictly followed before it is too late	1-9
	There is nothing wrong with premarital sexual intercourse (reversed)	1-9
Views on social inequality	Our society does NOT need tougher government and stricter laws (reversed)	1-9
	The facts on crime and the recent public disorders show we have to crack down harder on troublemakers, if we are going to preserve law and order	1-9
	I would like to bring the people above me on the ladder down a peg or two	1-7
	I would like to bring the people below me on the ladder up a peg or two	1-7
Religious views	How religious are you?	1-5
	It is likely that God, or some other type of spiritual non-human entity, controls the events in the world	1-7

Supplementary Table S4

Proportions of correct answers to comprehension questions for all six economic games, split by country.

Game	United Kingdom	United States
Game A (AI)	0.96	0.94
Game B (Equal)	0.95	0.93
Game C (Computer)	0.95	0.95
Game D (1:1 Fee-Fine)	0.95	0.94
Game E (DI)	0.96	0.94
Game F (Third-Party)	0.95	0.94