

Who is worthy to be trusted?

The relationship between fairness, harm, and perceived trustworthiness

Mariia Kuz (SNR: 2080888)

Master Thesis

April 1, 2022

Study program: MSc Economic Psychology

Academic year: 2021-2022

Supervisor: Dr. Florian van Leeuwen

Second reader: Dr. Ivana Vranjes

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Abstract

What do we pay attention to when deciding whether to trust a person or not? Especially when it comes to a stranger? This study was focused on two research questions: the presence of social categorization of harm and fairness as two distinct moral domains and impact of harm and fairness on perceived trustworthiness. The memory confusion task was employed to a sample of 65 English-speaking participants. The results indicate that individuals actually categorize harm as a distinct moral domain, while fairness not. Still, fairness may be a potential moderator between harm and perceived trustworthiness. Despite its exploratory nature, this study offers some insight into economic psychology as it shed light on the factors that influence the choice of potential partner that is one of the main concerns of modern economists.

Keywords: moral domains, harm, fairness, harm, trustworthiness, memory confusion paradigm.

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Who is worthy to be trusted?

The relationship between fairness, harm, and perceived trustworthiness

Recent studies dedicated to social trust show that the perception of an individual as more trustworthy is closely connected with moral judgments this individual makes (Everett et.al, 2016). But in fact, there are a few different perspectives about moral foundations on which we base our moral judgments. This current paper examines the relationship between the two moral domains (namely, fairness and harm) and trustworthiness. It begins with highlighting the importance of studying trustworthiness for economic psychology, then goes on to emphasize the role of moral domains in shaping the trustworthiness of others, elucidating the heterogeneity of moral cues and explaining the necessity of focusing exactly on harm and fairness. The last part focuses on the methodology of this study and the prospective plan of analysis.

On the reasons to study trust within economic psychology perspective

Modern social interactions are rooted in the economic environment, and this bond provokes an enormous number of social dilemmas (Dawes, 1980; Manski, 2000). In this regard, cooperation is considered as one of the most important aspects that leads to a well-functioning economy (Ostrom, 2010). In turn, cooperation is tightly connected to the concept of trust or its absence (Bauer et.al, 2019; Bouma et.al, 2008). This can relate to both everyday economic problems (for example, should you give a friend a loan or not?) and issues of trust on a more global level (is it worth working with a particular organization? Will the management of this company deceive or cheat on you?). Thus, in the new global economy trustworthiness has become a central issue for establishing human relationships (Cook & State, 2017). While economists are concerned with the practical consequences of (dis)trust (Kalish et.al, 2021), social psychologists are more concentrated on the factors that influence the choice of who is worthy to be trusted, such as facial impressions (Jaeger, et.al, 2020;

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Todorov et al, 2009), neurobiological factors (Baumgartner et.al., 2008), ethnicity features (Birkás et.al, 2014), personality traits (Colquitt & Salam, 2015) and moral domains (Haidt, 2007). The latter is the focus of the present study.

Trust in the perspective of moral domains

On the contrary of different moral cues

When it comes to social categorization, morality seems to be one of the most important keys to understanding why someone is being perceived as good or bad (van Leeuwen & Penton-Voak, 2012; Wojciszke et al., 1998). Nevertheless, to date, there has been little agreement on what exactly morality domains are and what kinds of moral cues shape our perceptions.

According to one of the most prominent theories in this area, Haidt's Moral Foundation Theory, there are six main foundations of morality: care/harm, fairness/cheating, loyalty/betrayal, authority/subversion, sanctity/degradation, and liberty/oppression (Haidt, 2012; Haidt & Graham, 2007). It implies, that when it comes to making a moral judgment, some individuals take into account one of these moral concerns (for instance, "if a person hit someone, I would rather think that they are bad" – judgment, connected with care/harm domain). The proposed theory has much in common with the earlier ones (Graham et.al, 2009). For instance, the domains of care/harm and fairness/cheating are pretty similar to Shweder et al.'s (1997) ethics of autonomy. The theme of harm and fairness was also raised by Turiel (1983) with his moral domains: according to his theory, morality includes concepts of physical and psychological harm, as well as fair distribution of resources, freedoms, and rights.

However, there are some theories that consider only fairness or only harm as the basis for moral judgments. Thus, the morality-as-cooperation theory avoids the harm dimension (Curry, 2019), while the theory of dyadic morality overlooks fairness (Schein & Gray, 2017).

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Simultaneously, empirical evidence provides the information on the importance of both issues of harm and fairness, but still tend to consider only one of them: trolley dilemmas concern primarily the questions of harm (“who should be saved and who should we sacrifice?”) (Greene et al., 2009), whereas economic games like the Ultimatum Game concern fairness against self-interest (Fehr & Gächter, 2002; Vavra, Chang & Sanfey, 2018).

So, in different theories and empirical studies morality is usually operationalized as helping (vs. harming) or as playing fair (vs. cheating). But there is still some misleading in current literature: the theory above emphasizes the presence of these categories in different perspectives, but rarely considers it in conjunction. In this way, it still raises a question on the necessity to distinguish between harm and fairness. Therefore, the first research question of the current study is as follows: do individuals themselves distinguish harm and fairness as two independent moral domains? Then, taking into account empirical evidence approving that while solving moral dilemmas people refer to harm (“The Trolley Problem”) or fairness (“the Ultimatum Game”) we hypothesize that:

H₁: People perceive harm and fairness as two distinct moral domains.

How are fairness and harm related to trust?

Overall, people tend to trust people who seem to be more moral than not (Haidt, 2012). However, we have already defined, that moral dimensions are highly heterogeneous depending on the paradigm we prefer, and the main difference between various theoretical perspectives is the presence of harm or fairness. So, what if one person is moral in relation to harm, but violates the principles of fairness (and vice versa)? The remaining question is who do we trust more? To those who do not harm? Or to those who are fair?

Although fairness is more frequently considered as a crucial moral cue, very little is currently known about its’ impact on social trust. The existing body of research on trust and fairness has only focused on the impact of trust on perceived fairness but not vice versa

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

(Bianchi, et.al, 2014). However, the impact of harm on trust is more well established from previous studies (Siegel, et.al, 2019). Thus, extensive research by Paul Bloom (2013) has shown that even 3-month-olds are sensitive to the harm domain: they have more trust in subjects who do not harm rather than those who behave violently. Moreover, the study by Everett et.al. (2016) proposes that we trust people who use the deontological argumentation more than those who prefer the consequentialist one. But the most interesting part of their research is that they used dilemmas that concern only the harm dimension. So, they conceptualized that we trust people who ignore utilitarian consequences and reject to kill anyone (deontologist) more than those who prefer to kill one in order to save five (consequentialism). In this regard, the second hypothesis was proposed:

H₂: Harmless individuals are perceived as more trustworthy compared to fair individuals.

Method

Design

Overall agenda

The research was conducted in a 2x2 within-subject experimental design (memory confusion task) with two variables (fairness and harm) and two levels for both of them (fair/unfair and harmful/harmless). Also, recall errors, the trustworthiness of each target, and sociodemographic data were measured.

Memory confusion task

To understand whether participants take into consideration harm and fairness as distinct moral domains there was used the *memory confusion paradigm* (van Leeuwen et al, 2012). In general, the task consists of the presentation and recall phases. During the presentation phase, participants are asked to form an impression of the target individuals. The case is that pictures of individuals are shown along with statements representing studied

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

categories (like harm and fairness). When the second step begins, participants need to remember “who-said-what” (Taylor et.al, 1978): in other words, they need to link a target with the statement.

Then, the researcher calculates two types of mistakes: within-category and between-category ones. Within-category errors occur when a participant attributes a statement to the wrong target but this target “said” the sentence in the same moral domain (for instance, if we consider harm category, then “correct” target “said” something harmful, but the participant chose another target, who also “said” something harmful). Between-category errors mean that participant attributes a statement to a target from another moral domain (if the “correct” target said something harmful, but the participant chose the target who “said” harmless statement). So, if there are significantly more within-category errors than between-category errors then categorization is based on the particular moral domain.

Materials and procedure

The study was conducted online on the Qualtrics platform after approval by the Ethics Review Board.

First, participants were asked to fill informed consent form containing information about the purpose and procedure of the study. Next, to start the survey, they had to pass the selection for sociodemographic characteristics (fluency in English, country of residence), and only afterwards the main survey was started.

Then the presentation phase began. Participants were shown eight targets with captions related to fairness or harm. These statements were developed based on the Moral Violations Vignettes (Clifford et al., 2015) and Moral Judgement Items and Taboo Trade-Off Items (Graham et al., 2009). The targets were taken from the Psychological Image Collection at Stirling (<http://pics.psych.stir.ac.uk/>) and consist of neutral facial expression photographs of White males. The same targets were used in the Kharitonenko’s paper (2021). But in our

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

study, we counterbalanced the statements and targets: for instance, if in Kharitonenko's study a target was harmful and fair, then we attributed to it a harmless and cheating judgment and vice versa (see Appendix for all the targets and statements). Then, after the filler task (participants will be asked to remember countries from European Union), participants had to link the target with the statement (*recall phase*). Afterward, *recall errors* (within-category and between-category) were calculated and used as a dependent variable in distinguishing the dimensions (fair and harm). Then, to correct for a higher probability of making a between-category rather than within-category error, between-category errors were multiplied by 0.75.

In addition, respondents were asked to evaluate the *trustworthiness of each target* to indentificate what influences trust more strongly: harm or fairness (if any). Trustworthiness was measured by the question “How trustworthy is this person?” with a 7-point Likert- scale from 1= extremely untrustworthy” to “7 = extremely trustworthy” (Everett et al, 2016).

At the end of the survey, there was a debriefing statement with the researchers' contact details (in case participants have questions about the study).

Participants

To calculate the sample size G*Power version 3.1.9.7 was used. An ‘A Priori’ analysis for paired samples t-tests showed that the sample must consist of at least 147 participants both sexes to achieve the power of .95 and to detect an effect size $d_z = .3$ at $\alpha = .05$. Participants were controlled by age (all participants must be over 18) and fluency in English (self-report). Respondents were recruited via social networks (non-probability convenience and snowball sampling).

A total of 124 responses were received, 59 of which were deleted as outliers, incomplete responses, or responses with failed attention check. Therefore, the final sample consisted of 65 respondents aged 18 to 66 years ($M = 26.41$, $SD = 6.19$), 42 females, 22 males, 1 participant chose not to reveal their gender. Most of the respondents rated their

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

English skills as “fluent” ($N=47$), 9 individuals stated that they have average skills, 9 mentioned that they are native speakers.

Analysis strategy

The data was analyzed via SPSS software. Before starting the main analysis, the data was checked for abnormalities and outliers.

To test a hypothesis about the distinction between harm and fairness domains there was conducted a *Wilcoxon Signed Ranks Test* to compare the within-category (fairness or harm) and between-category (fairness or harm) errors. The within-category error means that the participant makes a mistake in the sense of choosing a target that did not ‘say’ the asked statement, but at the same time, this target ‘said’ a statement that belongs to the same dimension of morality. For instance, when a participant sees a statement about harm but matches it with the person who said another statement, but still about harm. Between-category errors occur when the participant instead of the person who spoke about harm chooses the one who spoke about fairness. If categorization exists in a fair or harmful dimension, participants are expected to make more within-category than between-category errors for this dimension (van Leeuwen et al, 2012).

To test the second hypothesis about the impact of each dimension on trustworthiness a *2 × 2 repeated-measures ANOVA* was performed (two IVs, each with two levels: harmful/harmless and fair/cheater).

Results

Preliminary analysis

Before proceeding to the main analysis, the key dependent variables were tested for normality of distribution by the Kolmogorov-Smirnov test. Due to the fact that the real distribution of the errors-variables (both within and between harm/fair categories) deviated from the normal one ($p < .05$), to check the first hypothesis about the distinction between

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

harm and fairness moral domains non-parametric tests should be used (Wilcoxon Signed Ranks Test instead of Paired Samples T-Test). As for the second hypothesis variables (evaluations of targets' trustworthiness), the distribution is normal ($p > .05$), so we proceeded with 2×2 repeated-measures ANOVA as it was planned before.

Correlation analyses showed that age was not related to the main study variables, min $r = -.08$, max $r = .21$, $p > .05$. Gender also turned out to be non-significantly correlated with the main variables, r ranged from $-.02$ to $.21$, $p > .05$ (see Table 1 in the Appendix).

Main analysis

Harm and fairness as distinct moral domains

The Wilcoxon Signed Ranks showed that we can partially support the Hypothesis 1. It was found that there is indeed a harm categorization as far as the number of within-harm errors (*Mean Rank* = 32.08, *Sum of Ranks* = 1283) was higher than the number of between-harm errors (*Mean Rank* = 30.45, *Sum of Ranks* = 670), $z = -2.15$, $p = .03 < .05$. Nevertheless, our results suggest that there is no categorization based on fairness: although the Mean Rank for within-fairness errors was higher (*Mean Rank* = 35.88, *Sum of Ranks* = 1004.5) than the number of between-fairness errors (*Mean Rank* = 25.8, *Sum of Ranks* = 825.5), the effect was not significant ($z = -.66$, $p = .51 > .05$).

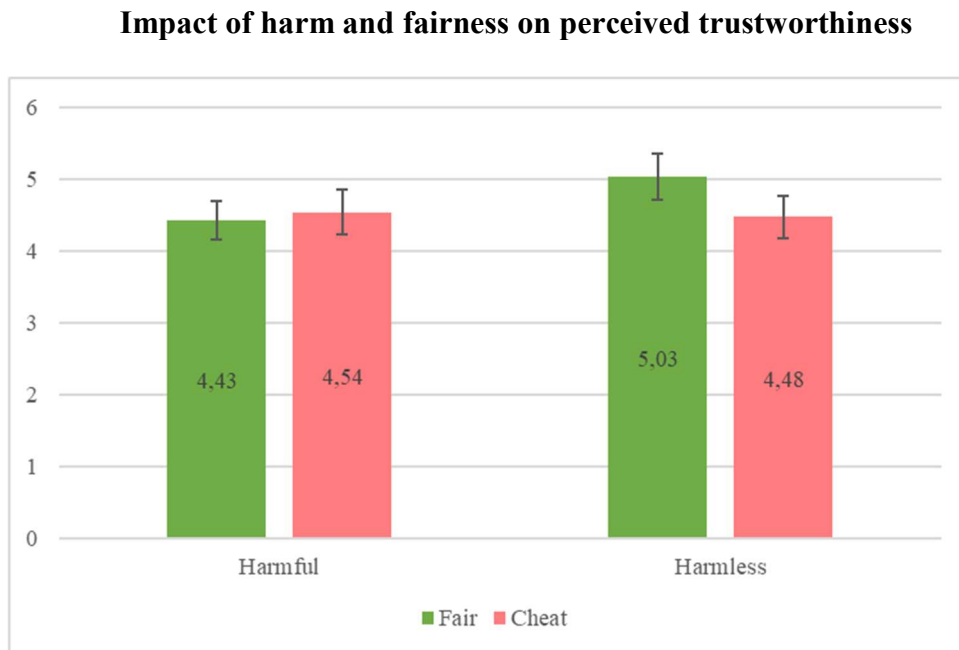
Impact of harm and fairness on perceived trustworthiness

The results show that there was no significant main effect of harm on perceived trustworthiness ($F(1,64) = 3.18$, $p = .08$, $\eta_p^2 = .05$). Approximately the same result was observed for fairness (boundary significance): $F(1,64) = 4.09$, $p = .047$, $\eta_p^2 = .06$. At the same time, it is most interesting to look at the interaction of these two factors ($F(1,64) = 10.33$, $p = .002$, $\eta_p^2 = .14$). The results below (see Figure 1) show that if the target was harmful, there is no matter if it is fair or not ($M=4.43$ for fair and $M=4.54$ for unfair). But if the target was harmless, then the evaluation depends on fairness: higher for fair ($M=5.03$), low for unfair

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

($M=4.48$).

Figure 1: Perceived trustworthiness mean scores



Discussion

This study was focused on two research questions: the presence of social categorization of harm and fairness as two distinct moral domains and impact of harm and fairness on perceived trustworthiness. So, the following conclusions can be drawn from the present research: (1) people actually categorize harm as a distinct moral domain, while fairness not; (2) trustworthiness of a person does not depend on harm and fairness separately: individuals equally evaluate both harmful/fair targets and harmless/unfair ones. However, if interaction of factors is taken into account, trustworthiness is higher in the most “moral” individuals (harmless/fair).

The evidence from this study support the idea of dyadic morality theory based on harm as a main moral dimension (Schein & Gray, 2017) and calls into question the existence of fairness as a distinct moral cue like it was proposed in Haidt’s Moral Foundation Theory (2012). Nevertheless, when it comes to assesment of trustworthiness, interaction between harm and fairness plays the role: if individual is harmless, then their perceived

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

trustworthiness is higher when they are also fair. In line with previous research (Everett et al., 2016), these findings indicate that people tend to trust more moral partners. Moreover, in general, it seems that harm more clearly affects the moral character of a person and the willingness to trust them: if the target is harmful, then it doesn't matter fair they are or not – the trustworthiness will be low. This evidence indirectly confirms our second hypothesis (“harmless individuals are perceived as more trustworthy compared to fair individuals”), despite the fact that the main effect for harm was insignificant.

Therefore, the empirical findings in this paper provide a new understanding of fairness in moral categorization as a possible moderator in the relations of different moral domains (future research needs to check the moderation role of fairness as far as it may strengthen the relationship between harm and perceived trustworthiness).

Finally, a number of important limitations need to be considered. First, the relatively small sample size that was caused by a low completion rate. Possible reasons for this may lie in the rather long questionnaire (most of respondents ended the survey during the recall phase). Furthermore, many respondents were non-native speakers, which constituted an additional cognitive load for them. Secondly, the main effect for fairness appeared in the expected direction (within-fairness errors > between-fairness errors). Thus, it is possible that due to the small sample size the statistical power was not enough to detect a significant result while in general the effect exists.

Notwithstanding these limitations, the study offers some insight into economic psychology as it sheds light on the factors that influence the choice of potential partner that is one of the main concerns of modern economists.

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

References

- Bauer, P. C., Keusch, F., & Kreuter, F. (2019). Trust and cooperative behavior: Evidence from the realm of data-sharing. *PLOS ONE*, *14*(8), e0220115.
<https://doi.org/10.1371/journal.pone.0220115>
- Baumgartner, T., Heinrichs, M., Vonlanthen, A., Fischbacher, U., & Fehr, E. (2008). Oxytocin Shapes the Neural Circuitry of Trust and Trust Adaptation in Humans. *Neuron*, *58*(4), 639–650.
<https://doi.org/10.1016/j.neuron.2008.04.009>
- Bianchi, E. C., Brockner, J., Van den Bos, K., Seifert, M., Moon, H., Van Dijke, M., & De Cremer, D. (2014). Trust in Decision-Making Authorities Dictates the Form of the Interactive Relationship Between Outcome Fairness and Procedural Fairness. *Personality and Social Psychology Bulletin*, *41*(1), 19–34.
<https://doi.org/10.1177/0146167214556237>
- Birkás, B., Dzhelyova, M., Lábadi, B., Bereczkei, T., & Perrett, D. I. (2014). Cross-cultural perception of trustworthiness: The effect of ethnicity features on evaluation of faces' observed trustworthiness across four samples. *Personality and Individual Differences*, *69*, 56–61.
<https://doi.org/10.1016/j.paid.2014.05.012>
- Bloom, P. (2013). *Just babies: The origins of good and evil*. Crown Publishers/Random House.
- Bouma, J., Bulte, E., & Van Soest, D. (2008). Trust and cooperation: Social capital and community resource management. *Journal of Environmental Economics and Management*, *56*(2), 155–166.
<https://doi.org/10.1016/j.jeem.2008.03.004>

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

- Clifford, S., Iyengar, V., Cabeza, R., & Sinnott-Armstrong, W. (2015). Moral foundations vignettes: A standardized stimulus database of scenarios based on moral foundations theory. *Behavior Research Methods*, *47*(4), 1178–1198.
<https://doi.org/10.3758/s13428-014-0551-2>
- Colquitt, J. A., & Salam, S. C. (2015). Foster Trust through Ability, Benevolence, and Integrity. *Handbook of Principles of Organizational Behavior*, 389–404.
<https://doi.org/10.1002/9781119206422.ch21>
- Cook, K. S., & State, B. (2017). Trust and social dilemmas: A selected review of evidence and applications. In P. A. M. Van Lange, B. Rockenbach, & T. Yamagishi (Eds.), *Trust in social dilemmas* (pp. 9–30). Oxford University Press.
<https://doi.org/10.1093/oso/9780190630782.003.0002>
- Curry, O. S., Jones Chesters, M., & Van Lissa, C. J. (2019). Mapping morality with a compass: Testing the theory of “morality-as-cooperation” with a new questionnaire. *Journal of Research in Personality*, *78*, 106–124.
<https://doi.org/10.1016/j.jrp.2018.10.008>
- Dawes, R. M. (1980). Social Dilemmas. *Annual Review of Psychology*, *31*(1), 169–193.
<https://doi.org/10.1146/annurev.ps.31.020180.001125>
- Essink, W. (2021). *Influence of fairness and reciprocity on impression formation and social trust* [MSc Thesis].
- Everett, J. A. C., Pizarro, D. A., & Crockett, M. J. (2016). Inference of trustworthiness from intuitive moral judgments. *Journal of Experimental Psychology: General*, *145*(6), 772–787. <https://doi.org/10.1037/xge0000165>
- Fehr, E., & Gächter, S. (2002). Altruistic punishment in humans. *Nature*, *415*(6868), 137–140. <https://doi.org/10.1038/415137a>

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

- Graham, J., Haidt, J., & Nosek, B. A. (2009). Liberals and conservatives rely on different sets of moral foundations. *Journal of Personality and Social Psychology*, *96*(5), 1029–1046.
<https://doi.org/10.1037/a0015141>
- Greene, J. D., Cushman, F. A., Stewart, L. E., Lowenberg, K., Nystrom, L. E., & Cohen, J. D. (2009). Pushing moral buttons: The interaction between personal force and intention in moral judgment. *Cognition*, *111*(3), 364–371.
<https://doi.org/10.1016/j.cognition.2009.02.001>
- Haidt, J. (2007). The new synthesis in moral psychology. *Science*, *316*(5827), 998–1002.
<https://doi.org/10.1126/science.1137651>
- Haidt, J. (2012). *The righteous mind: Why good people are divided by politics and religion*. Pantheon Books.
- Haidt, J., & Graham, J. (2007). When morality opposes justice: Conservatives have moral intuitions that liberals may not recognize. *Social Justice Research*, *20*, 98–116.
<https://doi.org/10.1007/s11211-007-0034-z>
- Jaeger, B., Oud, B., Williams, T., Krumhuber, E., Fehr, E., & Engelmann, J. (2020). *Trustworthiness detection from faces: Does reliance on facial impressions pay off?*
- Kalish, I., Wolf, M., & Holdowsky, J. (2021, 25 oktober). *The link between trust and economic prosperity*. Deloitte Insights. Geraadpleegd op 3 december 2021, van <https://www2.deloitte.com/us/en/insights/economy/connecting-trust-and-economic-growth.html>
- Kharitonenko, A. (2021). *To be fair or harmless? The relations of moral domains with cooperation* [MSc Thesis].

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

- Manski, C. F. (2000). Economic Analysis of Social Interactions. *Journal of Economic Perspectives*, 14(3), 115–136.
<https://doi.org/10.1257/jep.14.3.115>
- Ostrom, E. (2010). Beyond Markets and States: Polycentric Governance of Complex Economic Systems. *American Economic Review*, 100(3), 641–672.
<https://doi.org/10.1257/aer.100.3.641>
- Schein, C., & Gray, K. (2017). The Theory of Dyadic Morality: Reinventing Moral Judgment by Redefining Harm. *Personality and Social Psychology Review*, 22(1), 32–70.
<https://doi.org/10.1177/1088868317698288>
- Shweder, R. A., Much, N. C., Mahapatra, M., & Park, L. (1997). The “big three” of morality (autonomy, community, and divinity), and the “big three” explanations of suffering. In A. Brandt & P. Rozin (Eds.), *Morality and health* (pp. 119–169). New York: Routledge.
- Siegel, J. Z., Estrada, S., Crockett, M. J., & Baskin-Sommers, A. (2019). Exposure to violence affects the development of moral impressions and trust behavior in incarcerated males. *Nature Communications*, 10(1).
<https://doi.org/10.1038/s41467-019-09962-9>
- Taylor, S. E., Fiske, S. T., Etoff, N. L., & Ruderman, A. J. (1978). Categorical and contextual bases of person memory and stereotyping. *Journal of Personality and Social Psychology*, 36(7), 778–793. <https://doi.org/10.1037/0022-3514.36.7.778>
- The Enduring Relevance of Darwin’s Theory of Morality. (2013). *BioScience*, 63(7), 513–514. <https://doi.org/10.1525/bio.2013.63.7.2>
- Todorov, A., Pakrashi, M., & Oosterhof, N. N. (2009). Evaluating Faces on Trustworthiness After Minimal Time Exposure. *Social Cognition*, 27(6), 813–833.
<https://doi.org/10.1521/soco.2009.27.6.813>

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Turiel, E. (1983). *The development of social knowledge: Morality and convention*.

Cambridge, United Kingdom: Cambridge University Press.

van Leeuwen, F., Park, J. H., & Penton-Voak, I. S. (2012). Another fundamental social category? Spontaneous categorization of people who uphold or violate moral norms. *Journal of Experimental Social Psychology*, 48(6), 1385–1388.

<https://doi.org/10.1016/j.jesp.2012.06.004>

Vavra, P., Chang, L. J., & Sanfey, A. G. (2018). Expectations in the Ultimatum Game: Distinct Effects of Mean and Variance of Expected Offers. *Frontiers in Psychology*, 9. <https://doi.org/10.3389/fpsyg.2018.00992>

Wojciszke, B., Bazinska, R., & Jaworski, M. (1998). On the dominance of moral categories in impression formation. *Personality and Social Psychology Bulletin*, 24(12), 1251-1263.

[https://doi: 10.1177/01461672982412001](https://doi.org/10.1177/01461672982412001)

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Appendix

Stimuli: targets and statements

**Target 1**

- Harmful: It is okay to hurt others physically and emotionally.
- Unfair: My employees are working harder than ever, but instead of giving them a bonus, I'd rather keep the money for myself.

**Target 2**

- Harmful: It is okay to violate people's rights from time to time.
- Unfair: It is okay to take the credit for other people's hard work.

**Target 3**

- Harmful: It is okay to hit children for getting bad grades in school.
- Fair: I would feel uncomfortable cutting in a long line because it wouldn't be fair to those behind me.

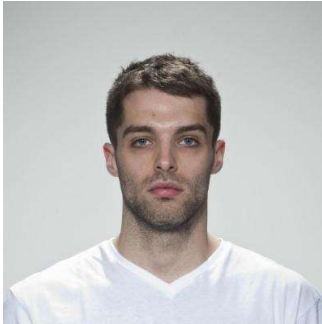
**Target 4**

- Harmful: It is okay to make cruel remarks to people about their appearance.
- Fair: All people deserve to be treated fairly.

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

**Target 5**

- Harmless: I can't stand cruelty.
- Unfair: It is okay to copy your classmate's work to get an A in class.

**Target 6**

- Harmless: Compassion for those who suffer is an important virtue.
- Unfair: It's okay to cheat others for your own benefit.

**Target 7**

- Harmless: People and especially children should be protected from harm at all costs.
- Fair: It is important that employees are paid appropriately according to their work.

**Target 8**

- Harmless: I am outraged when I see people hurting animals or children.
- Fair: When working on a project, I make sure to contribute as much as my teammates do.

RELATIONS OF FAIRNESS, HARM AND TRUSTWORTHINESS

Table 1*Descriptive Statistics and Correlations for Study Variables*

Variable	M	SD	1	2	3	4	5	6	7	8	9	10
1. Within-Harm Errors	2.68	1.29	-									
2. Within-Fair Errors	2.45	1.23	.12	-								
3. Between-Harm Errors	2.05	1.29	-.524**	.02	-							
4. Between-Fair Errors	2.23	1.12	-.16	-.341**	.341**	-						
5. Trust Harm&Fair	4.43	1.09	-.10	-.06	.07	.17	-					
6. Trust Harm&Cheat	4.54	1.29	-.08	-.22	.20	.269*	.632**	-				
7. Trust Not Harm&Fair	5.03	1.32	-.07	.08	-.265*	-.250*	.17	.17	-			
8. Trust Not Harm&Cheat	4.48	1.19	-.01	.05	-.11	-.01	.335**	.370**	.396**	-		
9. Age	26.42	6.20	-.08	.10	-.01	.11	-.05	.05	.11	.21	-	
10. Gender ^a	n.a.	n.a.	-.02	.00	.10	.14	.02	.21	.00	.21	.517**	-

Note. ^a1 = female, 2 = male; N=65; * $p < .05$, ** $p < .01$.