

Common-sense Disutility of the Erroneous Verdicts in Criminal Trials*

Kwangbai Park[†]

Yoori Seong

Chungbuk National University

As expressed by the famous Blackstone's ratio, the beyond reasonable doubt standard of proof is based on the law's primary motivation to avoid false conviction at the expense of increasing the probability of false acquittal. In contrast, jurors may have common-sense motivation to avoid both types of error. With 100 juror-eligible adults in Korea, the present study demonstrated that utilities of the two types of decision error were evaluated relative to those of the correct decisions rather than each other. The utility of false conviction was evaluated relative to that of a correct acquittal of an innocent defendant, and the utility of false acquittal in relation to that of a correct conviction of a guilty defendant. If this psychological configuration of the utilities is held by jurors in the courtroom, it suggests that they may have double standards for the fact-finding; one to decide on the question of guilt under the presumption of innocence and the other to decide on the question of innocence deduced from the presumption of guilt. Double standards will increase the frequency of punishing innocent defendants.

Key words : standard of proof, false conviction, false acquittal, utility

* This work was supported by the research grant of Chungbuk National University in 2014 and by the National Research Foundation of Korea Grant funded by the Korean Government (NRF-2013S1A5A2A03044871).

[†] Correspondence concerning this article should be addressed to Kwangbai Park, Department of Psychology, Chungbuk National University, Mt. 48 Gaeshin-dong, Cheongju City, Chungbuk, Republic of Korea. Email: kwangbai@chungbuk.ac.kr

Common-sense Disutility of the Erroneous Verdicts in Criminal Trials

The Presumption of Innocence is a “bedrock axiomatic and elementary principle” (*Coffin v. States*, 159 US 432, 453, 1895, Inre Winship, 397 U.S. 358, 1970에서 인용) of the modern criminal law. One of the legal apparatuses to put the principle in the day-to-day practice of law is the beyond reasonable doubt standard of proof for criminal cases. The prosecution must prove the defendant’s guilt beyond a reasonable doubt (e.g., Mueller & Kirkpatrick, 2009). The moral rationale or motivation of the standard was succinctly expressed by the Blackstone’s famous adage, “*It is better that ten guilty persons escape than that one innocent suffer.*” The gist of the adage compares two types of error on a negative utility (expected costs) dimension: a false conviction is expected to cause far greater costs to the society than would a false acquittal, and thus committing a false conviction is much more undesirable than committing a false acquittal. From this skewed ratio of relative utilities of the two types of error, the law’s motivation underlying the legal standard of proof can be inferred, as to minimize the probability of false convictions at the expense of a possible increase in the probability of false acquittals (DeKay, 1996; Newman, 1993).

The Supreme Court of Korea has explained the beyond a reasonable doubt standard for trial judges repeatedly over the last 40 years (1982DO263, 1997DO974, 2004DO2221). However, the standard is notoriously difficult for people to understand and explain (Brown, 2000; Cronan, 2002; Ellsworth & Reifman, 2000;

Garvey, Johnson, & Marcus, 2000; Laudan, 2003). More importantly, jurors appeared to have difficulty applying the standard of proof as the law intends (Stoffelmayr & Diamond, 2000). Research efforts to understand causes of the widespread difficulty associated with the comprehension, interpretation, and application of the standard have focused largely on the linguistics used for the judicial instructions on the standard and the words surrounding the concept of reasonable doubt (Elwork, Suggs, & Sales, 1981; Nagel & Neef, 1977; Young, Cameron, & Tinsley, 2001; Hastie, Penrod, & Pennington, 1983; Horowitz & Kirkpatrick, 1996).

A source of the difficulty may also be found in a more fundamental level, motivation. As Whitman (2008) points out, the historical motivation for the standard of proof was religious; its purpose was to protect the souls of the jury members, who, in 18th century Europe, feared heavenly damnation for condemning or punishing. In the same vein, another source of the difficulty may be due to the individuals’ common-sense motivation in legal fact-finding. According to a series of survey conducted with actual jurors in Florida (Strawn & Buchanan, 1976; Buchanan, Pryor, Taylor, & Strawn, 1978), Michigan (Reifman, Gusick, & Ellsworth, 1992), and Wyoming (Saxton, 1998) in the USA, a great portion (1/5 to 2/3) of the jurors believed that defendant’s innocence must be proved with evidence before being acquitted. Jurors’ motivation to avoid false acquittal may be as strong as their motivation to avoid false conviction. If so, jurors may have difficulty applying the standard of proof as the law intends because their intention in legal

decision-making is different from that of the law.

The present study was to examine the possibility that potential jurors may have motivation to avoid both false incrimination and false exoneration, as opposed to avoiding one type of error at the expense of committing another type of error. Utility, or expected costs and benefits, is defined as the willingness to pay (Samuelson & Nordhaus, 2009). Thus, pronounced utility can be an indication of underlying motivation. An individual's motivation can be measured with the relative utilities of outcome that are subjectively evaluated by the individual. If an individual expects the disutility or costs of a type of decision error to be great, she is likely to be motivated to avoid the type of error in fact-finding. The focus of the present study is on the dimensions underlying the individual's evaluation of the utilities of decision outcomes. It is difficult to find an empirical study showing how people of Korea would evaluate the utilities of possible decision outcomes relative to one another. In the evaluation, the two types of error may be pitted against each other on a single disutility dimension, as assumed by the Blackstone's ratio and the beyond reasonable doubt standard. Alternatively, they may fall on separate dimensions, which will be resulted if individuals evaluate the utilities of erroneous decisions relative to those of the correct ones rather than each other.

To measure individual's quantitative interpretation of the legal standard, the respondents were often asked to rate values of subjective utilities (expected costs and benefits) associated with wrongful, as well as correct, decisions (Fried, Kaplan, & Klein,

1975; Nagel, 1979). Those subjective utilities are put into a theoretical formula to indirectly derive a latent cut-off value on the probability of guilt. The theoretical formula is based on the signal detection theory that assumes the two types of decision error (false conviction and false acquittal) as distinctive outcomes. The cut-off threshold obtained with the method often appears outside the range of zero-to-unity, which suggests logical and/or psychometrical flaws (Dane, 1985). The result from the present study will shed some light on the anomalies of the particular method to measure the cut-off value corresponding to the legal standard. The present study will show that those outcomes are not psychologically distinguished from each other in their utilities.

Method

Individual difference scaling (Carroll & Chang, 1970; Carroll & Wish, 1974) — a type of multidimensional scaling — was applied to relative utilities of trial outcomes evaluated by 100 juror-eligible adults in Korea. Multidimensional scaling is a set of mathematical procedures developed to visually represent the epistemological or phenomenological structure of dissimilarities or comparative judgments for a set of stimuli. Like factor analysis, it can also be used to uncover dimensions on which the relative judgments of the stimuli are based. Individual difference scaling transforms measured dissimilarities into Euclidian distances in a k-dimensional space, and provides additional information regarding the individual's reliance (weights) on latent dimensions underlying

their judgments.

Participants

After discarding the data from 5 participants with incomplete or double responses, the responses of 95 juror-eligible adults (over 20 years of age) adults (49 men and 46 women) were analyzed for this study. The participants were sampled from a large panel of approximately one million Internet users who are registered as respondents for one of the major public-opinion research firms in Korea. Age distribution was 16.8%, 32.6%, 34.7%, and 15.8% respectively for the 20's, the 30's, the 40's, and over 50 years of age. Office worker (41.1%) was the most frequent occupation of the participants, followed by homemaker (12.6%), specialist (12.6%), service worker (10.5%), self-employed (7.4%), and other (15.8%). For the 70.5% of the participants, the level of education was higher than or equal to college graduate. The participants received a small amount of on-line credit for the participation in this study.

Measures and Procedure

Data was collected through an online Internet survey. The randomly selected candidates for participation were first contacted via email with introductory information about the study. Upon their agreement to participate, a temporary URL of the study site, which constantly changed, was forwarded to them. Precautions to secure the internal validity of the data included identification of the sampled respondents (e.g., matching IP address of their computers with the background

information registered on gender, age, education, etc.), limiting the number of participation (e.g., blocking repeated connections from the same IP address or by the same respondent), setting the time interval allowed for reading instructions that should not be too short nor too long (e.g., allowing a web-page to remain on the respondent's monitor screen for a certain interval of time), and confirming the veracity of responses at the end with the respondent.

It took the average person approximately 10 minutes to complete the survey. Participants first read instructions on the purpose of the survey, and how to respond to the survey items. Before responding, they were informed that as a potential juror serving for a criminal trial, they would rate relative costs and benefits of the four possible trial outcomes: correct acquittal, false conviction, correct conviction, and false acquittal. The participants were instructed to evaluate the relative utilities of those outcomes for their personal conscience, as well as for the society at large.

The participants rated the utilities of the four trial outcomes on a 21-point ordinal scale labeled with -10 for *very large cost*, $+10$ for *very large benefit*, and 0 in the middle. Each of the outcomes was described as a question: "What is the relative cost or benefit of convicting (acquitting) a truly innocent (guilty) defendant? Select a negative number on the left of the zero point to the extent that it causes costs. Select a positive number on the right of the zero point to the extent that it causes benefits. If its cost and benefit are balanced or cancel each other out, select zero." In order to encourage the relative judgments, all of the four trial outcomes and

corresponding utility scales appeared simultaneously on a single monitor screen in the spatial order of false conviction, false acquittal, correct conviction, and correct acquittal. The participants responded on the four scales, in any temporal order convenient for them, by clicking at the appropriate utility point on each of the four scales. Reselection of the utility points was allowed any time while the response screen was on the monitor.

Analysis

Six absolute differences among the four utilities rated by each participant were computed. Since the utilities were rated on a scale with the end points of -10 and 10 , the absolute difference between the negative utilities of the two false outcomes (false acquittal and false conviction), and that between the positive utilities of the two correct outcomes (correct acquittal and correct conviction) may normally range from 0 to 10 . Those between the rated utilities of a false outcome and a correct outcome can range from 0 to 20 . In order to make the scale ranges comparable, the absolute difference between the rated utilities of the two false outcomes and that of the two correct outcomes was each multiplied by 2 . With those differences, a 4×4 dissimilarity matrix was constructed for each participant. The dissimilarity matrix was symmetric around the main diagonal of zeros.

In the analysis, the sample was randomly divided into a test group of 50 participants and a cross-validation group of 45 participants. The stack of the dissimilarity matrices from each group was

subjected to the scaling algorithm developed by Carroll and Chang (1970). Since each of the four decision outcomes as stimuli is dichotomously varied on two attributes, decision (convict versus acquit) and outcome (correct versus false), their relative utility was likely to be evaluated along the two attributes. Accordingly, a two-dimensional space with its measurement level treated as ordinal was specified for the scaling algorithm.

Results

The means ($n = 95$) of the rated utilities were 4.66 ($SE = 0.38$, 95% CI [3.91 , 5.42]) for correct acquittal (CA), -6.95 ($SE = 0.28$, 95% CI [-7.50 , -6.40]) for false conviction (FC), 6.11 ($SE = 0.40$, 95% CI [5.33 , 6.89]) for correct conviction (CC), and -8.79 ($SE = 0.20$, 95% CI [-9.18 , -8.40]) for false acquittal (FA). The orthogonal contrast between the utility ratings of the correct and the false outcomes (CC & CA versus FC & FA), $F(1,94) = 799.62$, $p \leq .01$, $\eta^2 = .90$, the contrast between the two correct outcomes (CC versus CA), $F(1,94) = 14.45$, $p \leq .01$, $\eta^2 = .13$, and the contrast between the two false outcomes (FC versus FA), $F(1,94) = 40.92$, $p \leq .01$, $\eta^2 = .30$, were all highly significant.

Individual difference scaling on the dissimilarity matrices of the test group and the cross-validation group yielded almost identical configurations of the four decision outcomes. The coordinates of the test (cross-validation) group on the first dimension were 0.96 (0.94), -0.89 (-0.92), 1.04 (1.06), and -1.10 (-1.07) for CA, FC, CC, and FA respectively. The coordinates of the test

(cross-validation) group on the second dimension were -1.40 (-1.40), -0.01 (-0.00), 1.43 (1.43), and -0.02 (-0.02) for CA, FC, CC, and FA respectively. With the two groups combined, individual difference scaling on the 95 dissimilarity matrices of utility yielded a two-dimensional configuration of the four decision outcomes as shown in Figure 1. Kruskal's Stress Formula 1 value (Kruskal & Wish, 1978) of the configuration was 0.04 indicating a good fit between the dissimilarities and the configuration (see Kruskal, 1964 for the rule of thumb criteria of the stress value. Stress value under 0.05 indicates "good" to "excellent" fit). For another fit index, the squared correlation between the dissimilarities in the data and the distances in the configuration was $.998$, indicating that the ordering of the utility dissimilarities were almost perfectly represented by the distances in the two-dimensional configuration.

The mean of the normalized weights of the sample ($n = 95$) was $.88$ on the first dimension and $.23$ on the second dimension (The normalized

weight can range from zero to unity on each dimension). Analysis of angular variation (Mardia, 1972; Mardia et al., 1979) revealed no significant gender difference in the individual weights, $F(1,93) = 3.01$, $p > .05$, $\eta^2 = .03$.

Discussion

The present study was an attempt to capture a snapshot (Figure 1) of the common-sense epistemological structure involving the generic concepts of trial outcomes that the potential jurors may bring to the courtroom without any prior knowledge about the specific characteristics of the case. With the purpose in mind, experimental manipulation and control was kept at the minimum so that the resulting snapshot would have an external validity to some degree. To rate the utilities of trial outcomes, participants may have had in their mind different typical cases, and different imagined cases may have inflated the

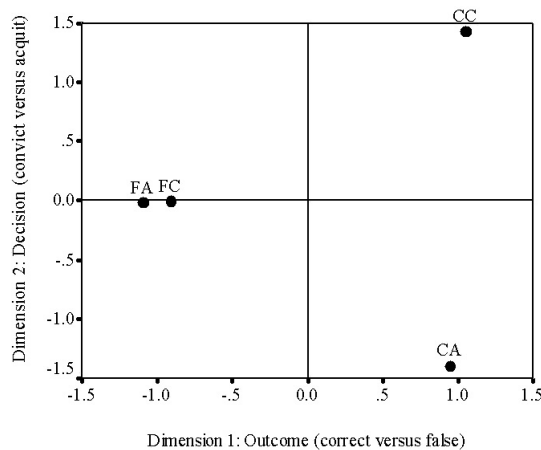


Figure 1. Configuration of the rated utilities of four decision outcomes ($n = 95$). FA = False Acquittal; FC=False Conviction; CA=Correct Acquittal; CC=Correct Conviction.

variance of random measurement error in the ratings. However, the error variance in the ratings was small, as suggested by the contrasts between the decision outcomes. The contrast between the two erroneous outcomes and the contrast between the two correct outcomes were, in spite of the small mean differences, highly significant because the standard errors in the ratings were small. With the modest sample size of $n=95$, the standard errors were only a fraction of one point on a 21-point scale. The small standard errors in the ratings indicate that the error variance caused by different imagined cases may be insignificant, if not nonexistent. People tend to imagine the most dangerous or heinous criminals known to them when they rate the cost of a false acquittal. And people tend to imagine the most innocent defendants they can imagine when they rate the cost of a false conviction. The most heinous and innocent cases imagined readily by different people would not vary greatly in nature in this modern era dominated by powerful media.

The two-dimensional configuration of the rated utilities of four decision outcomes reveals a psychological structure of decision outcomes that may influence laypeople's motivation underlying the common-sense standards of decision. The first dimension of Figure 1 can be clearly interpreted as the outcome (correct versus false) dimension. Less clearly, the second dimension seems to be the decision (convict versus acquit) dimension.

The rated utility of false acquittal (-8.79) was significantly more negative than that of false conviction (-6.95). If the negative utility of a type of error gives rise to an aversion to the type of error in decision-making, Korean jurors would

be particularly cautious to make the decision to acquit the defendant in court. However, the scaling of the utility ratings revealed that the discrimination between the two types of errors were relatively small, compared to the discrimination between the two types of correct decisions. The two types of errors were somewhat discriminated from each other on the outcome (1st) dimension that the participants relied heavily upon in their utility rating. The distance between the errors on the outcome dimension reflected the significant difference in utility rating between the two types of errors, as if a false acquittal is more false or erroneous than a false conviction. However, the two types of errors were almost indistinguishable from each other on the decision (2nd) dimension in spite of the fact that they are obviously different decisions. The lack of discrimination between the two types of error on the decision dimension is contrasted with the marked distance between the two correct decisions on the same dimension (Figure 1). Thus, the clustering of the two erroneous decisions on the decision dimension is not an artifact produced by the limited range of the scale used to measure the utilities.

When two different decisions, a conviction and an acquittal, are both erroneous, it may be difficult to weigh the negative utilities of them comparatively, probably because the variety of potentially negative effects of the errors on the involved individuals, the society, and social justice cannot be fully appreciated by common sense. This difficulty with which the negative utilities of the two types of erroneous decisions are compared with each other would, in turn, make it difficult

for individuals to appreciate the law's motivation that is often expressed by the Blackstone's famous adage.

Some research has attempted to derive laypeople's quantitative (i.e., probabilistic) interpretation of the beyond reasonable doubt standard of proof based on their judgments of outcome utilities (Fried, Kaplan, & Klein, 1975; Nagel, 1979). For the "statistical decision theory" method to measure individual's quantitative interpretation of the legal standard, the respondents were asked to rate values of subjective utilities (expected costs and benefits) associated with wrongful, as well as correct, decisions. Those subjective utilities are put into a theoretical formula to indirectly derive a latent cut-off value on the probability of guilt. The cut-off threshold has typically appeared near the chance level, and the standard for criminal trials was not meaningfully distinguished from the decision standards for civil trials (Arkes & Mellers, 2002; Connolly, 1987). Furthermore, the cut-off threshold obtained with the statistical decision theory often appears outside the range of zero-to-unity, which suggests logical and/or psychometrical flaws (Dane, 1985). The result from the present study sheds some light on the anomalies of the particular method to measure the cut-off value corresponding to the legal standard. The theoretical formula used to derive the latent cut-off value from rated utilities is based on the signal detection theory that assumes the two types of decision error (false conviction and false acquittal) as distinctive outcomes. The present study, however, shows that those outcomes are not psychologically distinguished from each other in their utilities.

The most prominent feature of Figure 1 is that the utilities of the two types of erroneous decisions were discriminated in comparison to those of the correct decisions, but not from each other. The utility of false conviction was evaluated relative to that of a correct acquittal of an innocent defendant, and the utility of false acquittal in relation to that of a correct conviction of a guilty defendant. If this psychological configuration of the utilities is held by jurors in the courtroom, it suggests that they may have double standards for the fact-finding. Jurors would first need a stringent standard not to convict the defendant erroneously. This standard would be to decide on the question of guilt (guilty or not guilty) under the presumption of innocence. At the same time, the juror would also need another equally stringent standard not to acquit the defendant erroneously. The second standard would be to decide on the question of innocence deduced from the presumption of guilt.

The results of this study have a number of limitations in external validity. First, to determine the generalizability of the results, other studies are necessary whose participants represent different demographics. Aversion to false acquittal may be relatively stronger among laypeople with low levels of education. A somewhat similar limitation of this study may be that all the participants were citizens of a country with a short history of trial by jury. Similar studies in countries with experienced juror-eligible adults may observe different dimensions underlying the evaluation of utilities.

References

- Arkes, H. R., & Mellers, B. A. (2002). Do juries meet our expectations? *Law and Human Behavior*, 26, 625-639.
- Brown, D. K. (2000). Regulating decision effects of legally sufficient jury instructions. *South California Law Review*, 73, 1105.
- Buchanan, R. W., Pryor, B., Taylor, K. P., & Strawn, D. U. (1978). Legal communication: An investigation of juror comprehension of pattern instructions. *Communication Monographs*, 26, 31-35.
- Carroll, J. D., & Chang, J. J. (1970). Analysis of individual differences in multidimensional scaling via an N-way generalization of Eckart-Young decomposition. *Psychometrika*, 35, 283 - 319.
- Carroll, J. D., & Wish, M. (1974). Models and methods for three-way multidimensional scaling. D. H. Krantz, R. C. Atkinson, R. D. Luce, & P. Suppes. (eds.) *Contemporary Developments in Mathematical Psychology, Vol.II: Measurement, Psychophysics, and Neural Information Processing*. San Francisco: W. H. Freeman and Company.
- Connolly, T. (1987). Decision theory, reasonable doubt, and utility of erroneous acquittals. *Law and Human Behavior*, 11, 101 - 112.
- Cronan, J. P. (2002). Is any of this making sense? Reflecting on guilty pleas to aid criminal juror comprehension. *American Criminal Law Review*, 39, 3, 1187-1259.
- Dane, F. C. (1985). In search of reasonable doubt. *Law and Human Behavior*, 9, 141 - 158.
- DeKay, M. L. (1996). The difference between Blackstone-like error ratios and probabilistic standards of proof. *Law and Social Inquiry*, 21, 95-132.
- Ellsworth, P. C., & Reifman, A. (2000). Juror comprehension and public policy: Perceived problems and proposed solutions. *Psychology, Public Policy & Law*, 6, 788 - 792.
- Elwork, A., Sales, B. D., & Suggs, D. (1981). The trial: A research review. In B. Sales (Eds.), *The Trial Process*. New York: Plenum.
- Fried, M., Kaplan, K. J., & Klein, K. W. (1975). Juror selection: An analysis of voir dire. In R. J. Simon (Ed.), *The juror system in America: A critical overview* (pp. 58 - 64). Beverly Hills, CA: Sage.
- Garvey, S. P., Johnson, S. L., & Marcus, P. (2000). Correcting deadly confusion: Responding to jury inquiries in capital cases. *Cornell Law Review*, 85, 628 - 633.
- Hastie, R., Penrod, S. D., & Pennington, N. (1983). *Inside the Jury*. Cambridge: Harvard University Press.
- Horowitz, I. A., & Kirkpatrick, L. C. (1996). A concept in search of definition: The effect of reasonable doubt instructions on certainty of guilt standards and jury verdicts. *Law and Human Behavior*, 20, 655-670.
- In re Winship*, 397 U.S. 358 (1970).
- Kruskal, J. B. (1964). Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29, 1 - 27.
- Kruskal, J. B., & M. Wish (1978). *Multidimensional Scaling*. Beverly Hills, CA: Sage.
- Laudan, L. (2003). Is reasonable doubt reasonable? *Legal Theory*, 9, 295-331.
- Mardia, K. V. (1972). *Statistics of Directional Data*. New York: Academic Press.
- Mardia, K. V., Kent, J. T., & Bibby, J. M. (1979). *Multivariate Analysis*. New York: Academic Press.
- Mueller, C. B., & Kirkpatrick, L. C. (2009). *Evidence* (4thed). Aspen: Wolters Kluwer.
- Nagel, S. S. (1979). Bringing the values of jurors in line with the law. *Judicature*, 63, 189 - 195.

- Nagel, S. S., & Neef, M. G. (1979). *Decision Theory and the Legal Process*. Lexington, MA: Lexington Books.
- Newman, J. O. (1993). Beyond “reasonable doubt.” *New York University Law Review*, 68, 979-1002.
- Reifman, A., Gusick, S. M., & Ellsworth, P. C. (1992). Real jurors' understanding of the law in real cases. *Law and Human Behavior*, 16, 539-554.
- Samuelson, P. A., & Nordhaus, W. D. (2009), *Economics: An Introductory Analysis* (19th Ed.) McGraw - Hill.
- Saxton, B. (1998). How well do jurors understand jury instructions? A field test using real juries and real trials in Wyoming. *Land and Water Law Review*, 33, 59-189.
- Stoffelmayr, E., & Diamond, S. S. (2000). The conflict between precision and flexibility in explaining “beyond a reasonable doubt.” *Psychology, Public Policy and Law*, 6, 769 - 787.
- Strawn, D. U., & Buchanan, R. W. (1976). Jury confusion: A threat to justice. *Judicature*, 59, 478-483.
- Whitman, J. Q. (2008). *The Origins of Reasonable Doubt: Theological Roots of the Criminal Trial*. New Haven: Yale University Press.
- Young, W., Cameron, N., & Tinsley, Y. (2001). *Juries in criminal trials* (Report 69). Wellington, New Zealand: New Zealand Law Commission.
- 1 차원교접수 : 2016. 01. 09.
수정원교접수 : 2016. 02. 23.
최종게재결정 : 2016. 02. 24.

