How Performance Pressure Influences AI-Assisted Decision Making

Anonymous ACL submission

Abstract

Many domains now employ AI-based decisionmaking aids, and although the potential for AI systems to assist with decision making is much 005 discussed, human-AI collaboration often underperforms due to factors such as (mis)trust in the AI system and beliefs about AI being incapable of completing subjective tasks. One potential tool for influencing human decision making is performance pressure, which hasn't been much studied in interaction with human-AI decision making. In this work, we examine how pres-012 sure and explainable AI (XAI) techniques interact with AI advice-taking behavior. Using an inherently low-stakes task (spam review clas-016 sification), we demonstrate effective and simple methods to apply pressure and influence 017 human AI advice-taking behavior by manipulating financial incentives and imposing time limits. Our results show complex interaction effects, with different combinations of pressure 021 and XAI techniques either improving or worsening AI advice taking behavior. We conclude 024 by discussing the implications of these interactions, strategies to effectively use pressure, and encourage future research to incorporate pressure analysis.

1 Introduction

028

034

042

With modern language models facilitating interaction with various AI systems, decision aids are now available across many industries (e.g., medical diagnoses Dilsizian and Siegel, 2014; Duron et al., 2021); financial management Zopounidis and Doumpos, 2002;and criminal recidivism risk McKay, 2020), and when used to complement human abilities, have the potential to outperform either the human or AI working alone. The potential is not necessarily realized, however, because of several challenges: debates on ethical resposibility of decisions (Smith, 2021; Busuioc, 2021; Johnson, 2021), the human ability to recognize when AI advice should be taken (Schemmer et al., 2023), mental models (biases) regarding AI performance and ability to perform well on subjective tasks (Clark et al., 2021; Jones-Jang and Park, 2023), and effects of how the AI advice is delivered (Steyvers and Kumar, 2023). 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

Many research directions thus aim to resolve these barriers to complementarity in human-AI performance, measured as *appropriate AI advice reliance*, which has two dimensions: taking correct AI advice and disregarding incorrect AI advice. Investigations include the effects of showing explanations from explainable AI (XAI) alongside AI system predictions (Bansal et al., 2021), introducing cognitive forcing functions when presenting AI advice (Buçinca et al., 2021), adjusting AI advice presentation methods (Rastogi et al., 2022), and adjusting task framing to account for biases about the types of tasks AI can work with (Castelo et al., 2019).

In AI-assisted decision making, the human makes the final decision, bearing full responsibility for its consequences. It is established that performance pressure from responsibility can influence decision making behavior (Ashton, 1990). But how does it influence *AI-assisted* decision making?

AI-assisted decision making experiments have considered tasks with stakes that are intrinsically high (loan defaults; Green and Chen, 2019) and low (speed dating; Castelo et al. (2019)), but the stakes have little tangible effect or implication for evaluators. Hence, we observe a gap in the literature of how people rely on AI assistants *under performance pressure*, that is, when stakes matter personally.

We believe this question is of special significance in the NLP research community, and not only in deployment scenarios. Modern NLP relies on eliciting high-quality data from humans to train models, often with systems in the loop. For example, Dynabench (Kiela et al., 2021) and ANLI (Nie et al., 2020) are datasets where humans work with AI models to create data through finding adversarial or interesting examples. Such datasets are often curated with low personal stakes, e.g., Wadhwa et al. (2024); Krishna et al. (2023, 2024); Lu et al. (2024) and Haduong et al. (2023) crowdsourced annotations and paid hourly wages. Could judiciously applied performance pressure influence the decisions of annotators building research datasets in ways that lead to improvements in data, and by extension AI evaluations and systems?

086

090

100

101

102

103

104

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

In this work, we seek to understand how performance pressure influences AI advice usage when the advice is provided as a second opinion. We recruit participants to decide whether a hotel review is genuine or deceptive and provide them with an AI advisor. We manipulate performance pressure in three different ways: by providing a bonus for correct answers, by deducting from the task compensation for incorrect answers, and by providing a bonus for correct answers within a time limit. We further investigate how performance pressure and different XAI techniques interact. Our results reveal a complex story. Under certain conditions, pressure can either improve or lower appropriate AI advice reliance, and XAI can sometimes mitigate negative effects of pressure.

Our contributions are:

- We demonstrate how to increase the stakes in an inherently low-stakes AI-assisted decision making setting; this approach can generalize to many pre-existing study designs.
- We show how XAI affects advice reliance (both positively and negatively), and how the effects interact with pressure, forming a complex picture about how human behavior changes. These findings suggest opportunities for designing adaptive decision-making environments when different XAI methods are available.
- We explore how pressure and confirmation bias can increase overreliance on AI advice and discuss implications of unintentionally encouraging people to trust AI *too* much.

2 Related Work

It is challenging to design systems that can assist with decision making, because people are influenced by many factors when taking advice, such as their personal expertise (Ronayne and Sgroi, 2019), the advisor's reputation (Yaniv and Kleinberger, 2000), or the style of advice delivery (e.g., inviting or broadcasting; Chhabra et al., 2013; Morrison et al., 2024), resulting in inconsistent advice taking behavior that can be challenging to predict. For example, even if advice is objectively high-quality (e.g., advice based on fact), it may still be discounted in the final decision (Wang and Du, 2018). AI advisors are expected to complement human decision making and result in higher collaborative performance, compared to individual performance, but recent studies have observed that over- and under-reliance on AI advice result in suboptimal collaboration (Bussone et al., 2015; Jacobs et al., 2021). 133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

Toward appropriate AI advice use. Algorithmic aversion has been shown to be task-dependent, in line with ideas about how well machines can perform on subjective tasks. When the task is subjective, e.g., predicting speed dating results, Castelo et al. (2019) found increased algorithmic aversion, as opposed to an objective task, e.g., predicting financial outcomes. Hypothesizing that people discount AI advice because they do not trust the AI system, researchers have used explainable AI (XAI) methods and shown the explanations alongside the AI advice. Many studies have observed XAI improving appropriate AI reliance (e.g., Panigutti et al., 2022; Ben David et al., 2021; Lee et al., 2023). Yet Fleiß et al. (2024) observed the opposite: when decisions were about quantifiable skills (e.g., work experience or command of English), rather than soft skills (e.g., diligence or ability to work in teams), adding explanations did not significantly increase AI advice reliance. Jiang et al. (2022) similarly observed XAI failing when the user is too uncertain. Another set of methods aim to mitigate poor AI reliance through cognitive forcing functions-interventions that cause a decision maker to engage in analytical thinking (Lambe et al., 2016). For example, Rastogi et al. (2022) successfully employ a cognitive forcing function to reduce anchoring bias-a bias where people weight earlier information higher-by adding a time delay before showing AI advice.

Decision making under pressure. An important environmental factor to consider is the influence of stressors on the human decision maker. Decision making often occurs under time stress or the weight of responsibility, for example. Different stressors can influence decision making in different ways (Henderson et al., 2024), and when multiple stres-

276

277

278

279

sors are present, their compound effect can present 184 itself in additive, synergistic, or antagonistic ways 185 (Hale et al., 2017). The influence of stress on AIassisted decision making is an understudied factor, 187 although in recent work, Swaroop et al. (2024) study how AI-assisted decision makers perform under time pressure, which emerges in real-world 190 settings like operating rooms and search and rescue 191 missions. They study when to provide AI advice in 192 an inherently high-stakes medical diagnosis task, 193 adapted to be approachable to laypeople. Gazit et al. (2023) studied AI-assisted decision making 195 under the pressure of responsibility and observed 196 how responsibility pressures overrode logical rea-197 soning, resulting in lower appropriate AI reliance. 198 The experiment setup involved surveying managers in business organizations, using experts with real responsibilities but asking about their behaviors rather than empirically observing them. Further work is needed to understand the role of responsibility and pressure in AI-assisted decision making.

Manipulating performance pressure. Performance pressure can be experimentally manipu-206 lated through different consequences, e.g., rewards 207 and reputation (Stoker et al., 2019). High-quality crowdworker data can be collected by using an appropriate financial incentive in the form of a fair 210 base pay and bonuses. A higher potential reward, 211 or bonus, can increase the pressure on the crowd-212 worker toward higher performance. A common 213 way of presenting the bonus is to frame it as a 214 215 gain, e.g., "if you do a good job, you can earn a bonus". Alternatively, the bonus could be framed 216 as a loss, e.g., "if you do a poor job, you will lose 217 your bonus". The literature in risk aversion, the 218 propensity to play it safe, and loss aversion, the 219 fear of losing out, has observed a stronger pressure effect from framing incentives as a loss rather than a gain (Merriman and Deckop, 2007). Grgić-222 Hlača et al. (2022) designed a study investigating 223 how trust in the AI advisor evolves and success-224 fully used the loss framing. In their experiment, users made AI-assisted decisions and updated their mental models of the AI behavior. Here, we are interested in studying the influence of external performance pressure stressors to encourage more appropriate AI reliance (i.e., correctly using AI advice to improve decision-making), rather than studying effects on user trust. 232

3 Experiment

From that extant literature, we form the following hypotheses:

H1: We can influence AI advice reliance by manipulating the environmental pressure. Increased performance pressure from monetary incentives framed as a loss will improve appropriate AI advice reliance, and increased pressure from time limits will reduce it.

H2. We can predict the influence of the performance pressure manipulation through measuring the risk aversion level of participants and their trust in the AI advisor. Participants with higher risk aversion will be more careful in their decision making. Participants with higher trust in the AI advisor will have more decisions aligned with the AI advice.

H3. Participants will spend more time considering the AI advice under higher pressure (both monetary and time), regardless of whether they change their decision, because they want to be more careful about their response.

H4. The positive effects from XAI will hold under pressure, potentially further increasing appropriate reliance of AI advice over no XAI.

To study how pressure influences AI advicetaking with and without AI explanation aids, we recruit Prolific¹ crowdworkers and task them with judging whether a hotel review is genuine or deceptive. We design a within-subjects experiment manipulating the task environment to apply pressure and run three experimental settings simultaneously, changing the availability of an explanation aid, allowing us to consider all of the above hypotheses at the same time.

3.1 Dataset

Of the many text annotation tasks available, we choose deceptive review classification because it has real-world importance, is not an inherently high-risk task (as compared to medical diagnosis), does not require expertise in the real-world setting (as compared to criminal recidivism), and likely has minimal relevance to our participants (e.g., the impact of predicting a review incorrectly has no personal effect on the participant). Hence, the pressure to perform well on this task must be primarily be external, necessitated by our experimental setup where we wish to simulate different levels of external performance pressure. It also

¹https://www.prolific.com/

371

372

325

parallels the annotation setting of data creation for NLP research.

We draw data from the Deceptive Opinion Spam Corpus (Ott et al., 2011, 2013), a binary classification task, which contains genuine hotel reviews from travel websites and deceptive reviews written by Amazon Mechanical Turk workers. The task is challenging: human performance is 55%—little better than random chance, ensuring that AI advice taking behavior we observe is not confounded by participants' prior knowledge or skill.

3.2 XAI Methods

281

282

290

291

296

301

304

307

312

313

314

315

318

319

320

323

324

We use two XAI methods: feature importance highlighting (LIME; Ribeiro et al., 2016), and natural language explanations produced by a generative AI (GenAI). LIME requires feature weights, thus we train an SVM classifier with tf-idf features. Our model achieves 86% accuracy on the test set using 5-fold cross-validation, in line with the SVM used by Schemmer et al. (2023). We do not disclose the model performance in our study to avoid user bias about objective performance metrics of the advice. For GenAI, we generated the explanation by prompting a large language model, ChatGPT², to explain why a review received a particular label (Appendix A). Note that this explanation is hallucinated, and the same review could receive a generated explanation for either label. We selected these approaches because LIME and GenAI are popular XAI methods used for text and studying AI-assisted decision making (Schemmer et al., 2023; Bansal et al., 2021).³

3.3 User Interaction

To measure the influence of AI advice under different stakes, we require a sequential decision making setup. For this reason, we use the judge-advisor system (JAS; Sniezek and Buckley, 1995). Under JAS, a user will first make a judgment alone, then receive advice, and finally make a second judgment (either confirming or adjusting their initial judgment). The sequential nature allows us to measure influence by comparing the final judgment with initial pre-advice judgments. Our interface design is heavily inspired by Schemmer et al. (2023) in order to establish our baseline conditions with previous work.

3.4 Independent Variables.

We measure demographics data (gender, education level, race), trust in AI and frequency of AI usage in work (5-point Likert), trust in the AI advisor (4 items), and risk aversion (Appendix A.2). Risk aversion is measured in two ways: the 10-item IPIP representation of the Tellegen (1995/2003) Multidimensional Personality Questionnaire⁴ (MPO; Appendix A.2) and the Holt and Laury (2002) Risk Assessment (HL; Appendix A.3.1) (10 items). MPQ asks subjects to rate their level of agreement with statements (e.g., "I avoid dangerous situations") using a 4-point Likert. HL contains a list of "gambles" where participants choose between "safe" and "risky" choices. Users are incentivized to answer truthfully on the survey to earn a bonus of up to 3.85USD.

3.5 Experiment Setup

Items. We sample 24 reviews from the test set, ensuring a balanced sample of genuine and deceptive reviews, and also of correct vs. incorrect AI predictions. All reviews had positive polarity. We select two additional reviews for practice: one where the AI is correct and one where it is incorrect. The reviews had a length of 45–120 words. Each pressure condition (details below) was assigned a balanced, random assignment of 8 reviews (2 of each {genuine, deceptive} \times {correct, incorrect}), and participants encountered pressure conditions in random order. We included two attention checks and rejected data from participants who failed both.

Subjects. A total of 302 participants were recruited on Prolific across three explainable AI (XAI) conditions. The recruitment conditions were 95% HIT acceptance rate, native English speaker, and limited to U.S. workers. After subjects accepted the task, they were directed to a consent form, completed a presurvey with demographics questions, questions about AI usage frequency and trust, and MPQ, then received instructions for the task. They completed two practice items and received feedback on the correctness of their decision to ensure they understood the JAS setup and also that the AI advice could be incorrect. For each item, reviewers decided whether a review was genuine or deceptive and rated their decision confidence

²Accessed January 20, 2025

³We have no hypotheses about different kinds of AI methods or their accuracy, nor about the faithfulness of the XAI explanations to the workings of the classifier. Hence we opted for a relatively simple but realistic classification system and widely-used XAI methods.

⁴https://ipip.ori.org/newMPQKey.htm

on a 7-point Likert scale, then received AI advice, 373 then were given the chance to update or confirm 374 their decision and confidence level. Participants did not receive feedback on the correctness of their judgments after the practice items, to ensure that trust in the AI system was held constant across items. After two practice tasks, which are excluded 379 from analysis, subjects judged all 24 reviews, then completed a postsurvey. The postsurvey contained questions to determine the level of trust in the AI advisor and HL. The study was approved by our institution's IRB, and participants were guaranteed a wage of 20USD/hr. Overall, the study took about 30 minutes per participant. The participants received 6USD base pay for completing the task and were aware of the bonuses. The average wage rate after bonuses was 33USD/hr. Participants whose performance resulted in underpayment received bonuses to meet the wage rate. 391

Experimental conditions. We use a withinsubjects design varying the type of pressure (baseline pressure, payment pressure, or time pressure). In the baseline pressure condition, participants are 395 informed they will receive a bonus of 0.5USD for every correct decision. In the payment pressure condition, participants are additionally informed they will lose 0.8USD for every incorrect answer. In the time pressure condition, participants must 400 make a correct decision within 30 seconds to re-401 ceive the 0.5USD bonus. A timer was displayed 402 to indicate remaining time. If participants ran out 403 of time, the timer would count down negatively. 404 We also use a between-subjects design to study 405 the effects of XAI decision aids (baseline, LIME, 406 or GenAI). The baseline subject group received 407 no explanation for the AI's prediction, the LIME 408 subject group received a LIME feature importance 409 explanation (where text is highlighted to indicate 410 its association with a label), and the GenAI sub-411 ject group received a natural language explanation 412 generated by ChatGPT (see §3.2). 413

> Appendix Fig. 3 shows an example of the interface with payment pressure and LIME explanations.

414

415

416

417

418

419

420

421

422

Dependent variables. The following measures are our dependent variables:

- Total accuracy of the final answers given by the human.
- Relative positive AI reliance (RAIR; Appendix A Eq. 1): the ratio of the number of

cases where the human relies on AI advice to correct their decision (i.e., they were incorrect before receiving advice and correct after) (Schemmer et al., 2023). 423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

• Relative positive self-reliance (RSR; Appendix A Eq. 2): the ratio of the number of cases where the human correctly maintains their judgment, disregarding the incorrect AI advice (Schemmer et al., 2023).

For all three measures, higher values are preferable.

4 Results

After filtering for failed attention checks, we collected responses from 99, 102, and 101 subjects across the three XAI conditions: baseline, LIME, and GenAI, respectively.

Demographics and surveys. Participants' demographics reported were: 51% male, 47% female, and 2% other; aged 18-30 (34%), 31-45 (46%), 46-60 (16%), and 61+ (4%); 70% were white, 16% were black, 6% were asian, and the remainder were mixed; and 31% completed a bachelor's degree, 28% completed high school, 13% had an associate's degree, and the remaining completed graduate degrees. The MPQ survey found 65% of respondents were risk-loving and 35% were risk-averse. The HL survey found that 70% of respondents were risk-averse, in direct contrast to the MPQ results. Participants generally trusted AI (55% of subjects rated ≥ 4 , 33% rated 3, 12% rated ≤ 2) and frequently used AI to help with their work (59% responded ≥ 4 , 18% rated 3, 23% rated ≤ 2).

Figure 2 summarizes accuracy after receiving AI advice across all conditions. We note that GenAI advice increases accuracy slightly; time pressure slightly decreases accuracy, mirroring Swaroop et al. (2024)'s limited findings of how time pressure influences AI advice-taking.

Pressure via monetary loss has varying effects on AI advice usage, and time pressure lowers appropriate AI advice usage (H1). Figure 1 summarizes the difference in relative positive AI- and self-reliance (respectively, RAIR on the left and RSR on the right), across conditions. We observe varied results with the payment pressure condition. When a natural language explanation is given, RAIR drops by 5% on average; RSR decreases in both no-XAI (6%) and GenAI conditions (8%).



Figure 1: GenAI improves appropriate AI advice reliance, but pressure has a predominantly negative effect.



Figure 2: Accuracy after receiving AI advice.

Time pressure has the strongest influence, decreasing RAIR and RSR for both XAI methods. Participants completed these tasks faster, averaging 20 seconds per task compared to the baseline and payment conditions (compared to 30 seconds), despite being given 30 seconds to complete each task. Hence we attribute the performance drop in the time pressure condition to rushing.

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

We build random effects regression models to investigate how pressure and XAI conditions influence observed RAIR and RSR differences (Appendix B). Neither variable was a predictor for RAIR, but XAI condition is a significant predictor for RSR (p < .05). Confirming our observations that time pressure lowers overall accuracy, the regression models found pressure condition to be a significant predictor (p < 0.05) for accuracy.

Overall, we find that payment pressure has an observable negative effect on appropriate AI advice reliance and overall accuracy, rejecting our hypothesis that performance pressure will improve appropriate AI reliance. The influence is most pronounced in the GenAI condition. The time pressure condition appears to have stressed participants out, influencing them to rely more on the AI advice, resulting in lower RAIR, RSR, and overall accuracy.

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

Risk aversion scores are not a predictor for AI advice usage, but AI advisor trust is (H2). We build random effects regression models to investigate how well the AI advisor trust and risk aversion surveys predict AI advice usage (Appendix B). Although the two risk aversion surveys had opposing results for participants (the MPQ personality survey rated participants as risk-loving whereas the Holt-Laury survey rated the same participants as risk-averse, and vice versa), neither risk score was a predictor for RAIR, RSR, or overall accuracy. We attribute this result to the minimal influence of payment pressure on the dependent variables. Participants' trust in the AI advisor are a significant predictor for both RAIR and RSR (p < .05), though general trust in AI was not a predictor. Neither risk nor trust measures were predictors for overall accuracy. From the postsurvey: 10 participants stated they relied on the AI after realizing the first few judgments aligned with the AI advice, indicating trust in the AI advisor, and 7 participants stated they guessed or followed their gut feeling.

Payment pressure increases time spent considering AI advice (H3). Although the influence of payment perssure does not have a strong influence on RAIR and RSR, we do observe different decision-making behavior. In the payment pressure condition, participants on average spent 40% longer making the initial decision and 10% longer considering AI advice before making the final decision, as compared to the baseline pressure con-

619

620

621

622

623

624

625

575

576

dition. Despite the increased time spent on tasks,
overall performance was largely unaffected, with a
minor increase when no XAI is present and a minor
decrease when GenAI is present (Figure 2).

532

533

538

539

540

541

543

544

545

548

549

551

554

XAI can mitigate negative effects of pressure (H4). In contrast to Schemmer et al. (2023), whose work found that LIME improved RAIR and had a minimal effect on RSR, our results show LIME decreases RAIR compared to no-XAI and reduces RSR in the time pressure condition. Yet LIME synergizes well with payment pressure, improving RAIR over the baseline and maintaining RSR. The RSR change aligns with two participants who stated they relied on the AI significantly and entirely during the time conditions. GenAI explanations improved RAIR and RSR in the baseline condition, but pressure canceled out this effect. However, GenAI marginally improved overall accuracy compared to the other two explanation conditions, indicating the potential for GenAI-style XAI in challenging tasks.

4.1 How did participants judge reviews?

Our postsurvey asked participants to describe how they determined whether a review was genuine or deceptive. Several subjects admitted to guessing because they could not determine any relevant features indicating review quality (genuine or deceptive). Twenty-four subjects checked for grammar, typos, and punctuation, and 23 examined the specifity of the reviews. Eight subjects "tried to feel human emotions", suggesting they associated deceptive reviews with being algorithmically generated, which describes algorithmic aversion for subjective tasks (Castelo et al., 2019). Two subjects stated they were determined not to "lose some bonus", indicating that the payment pressure condition affected their behavior.

5 Discussion

This work investigated how environmental pressure, combined with XAI, influences the way crowdworkers complete an AI-assisted decision making task. A large body of literature has focused on improving AI reliance isolated from environmental factors such as pressure, but AI-assisted decision making often occurs under pressure. We demonstrated two simple methods for inducing pressure through payment framed as a loss and limiting time.

5.1 AI advice reliance under pressure

There are many types of environmental pressures that an AI-assisted decision maker can be under, and interaction effects can vary widely. When looking at payment pressure, our findings show a subtle effect of how AI advice reliance changes. Without explanations, the pressure conditions had minimal influence on RAIR but negative influence on RSR. With LIME, RAIR improved and RSR remained the same with payment pressure, but RAIR and RSR both decreased with time pressure. With GenAI, RAIR and RSR consistently decreased.

The RAIR decrease in LIME in the baseline pressure condition compared to no-XAI suggests LIME decreased trust in the AI advice. It associated individual words with a prediction that may not have made sense (e.g., "is" is associated with being genuine; Figure 3b), but users reported associating typos, excessive punctuation, and grammar with their prediction (§4.1). Payment pressure appears to mitigate the negative effects of LIME, improving RAIR over the LIME baseline (but not over the no-XAI baseline) without changing RSR. This result points toward the potential for XAI techniques in the same style as LIME to help mitigate or further improve AI-assisted decision making behavior on tasks where humans are close to random.

GenAI has the worst interaction with pressure. RAIR decreases to the level of no-XAI in the payment condition and to the level of LIME in the time condition. We suspect that the payment pressure increased skepticism in AI advice, since natural language explanations can appear generic (e.g., "The review is deceptive because it overly praises the hotel without mentioning any potential downsides"). The time spent considering AI advice increased over the baseline, which can be interpreted as payment pressure influencing people to take more care in their final decision, or reducing their trust in the AI advice in a similar manner to how narcissism has a negative relationship with advice taking (Kausel et al., 2015; O'Reilly and Hall, 2021).

One XAI method can be more effective than another depending on the task or stakeholders (Jiang et al., 2022), and our results show that this is the case even under different pressure conditions. Furthermore, the negative effect of pressure can override the potential benefits of XAI, as seen in the GenAI RAIR decrease.

Time pressure has a largely consistent, negative impact on RAIR, RSR, and overall performance.

The behavior change mirrors the phenomenon of "choking," where a human experiences performance decline in critical high pressure settings, often attributed to anxiety. Participants were anxious about completing the tasks within the 30 second time limit, and rushed to complete them faster (averaging only 20 seconds, compared to 30 when working without time pressure). As a result, they relied heavily on the AI advice. One subject stated: "For the timed cases, I ended up relying on the AI's decision and just went with it. However, for the other parts, I looked out for typos and missing punctuations to detect human reviews."

626

627

632

637

639

645

647

651

666

670

671

672

Swaroop et al. (2024) observed that slowing down AI overreliers in higher time pressure environments could improve their performance, and since our participants had plenty of time to spare, we expect such a strategy would be helpful in this setting. Emotional regulation strategies can be effective in improving performance when anxiety is high (Balk et al., 2013), which can help alleviate anxiety over running out of time (and hence the possibility of earning a bonus). Methods from the distraction model may also help: if choking occurs due to cognitive overload (rather than anxiety), an attention shift could refocus the user to pay attention to relevant information and avoid choking (Hardy et al., 2001; Mullen et al., 2005; Eysenck, 2012; Nieuwenhuys and Oudejans, 2012). For example, adding a third color to LIME explanations, unrelated to the prediction, could be a distractor.

5.2 Practical Advice

A motivating use case for this work is AIassisted crowdsourced annotation and data elicitation projects in NLP research. In summary, we recommend that GenAI-based explanations be explored for benefits to appropriate AI advice reliance and accuracy. For projects interested in exploring how annotations or elicited data change under different pressure environments, e.g., users may prefer curt dialogue interactions in high pressure environments but elaborate responses otherwise, care should be taken to ensure confirmation bias does not increase AI trust throughout the course of the project. Performance pressure may be a useful tool, but we did not find consistent benefits.

6 Future Directions

673Our work has illustrated how environmental stres-674sors and XAI can alter AI-assisted decision making

behavior. The mixed effectiveness of incorporating XAI under pressures reinforces the findings of Jiang et al. (2022) and Swaroop et al. (2024) where AI advice or XAI are only effective for some people, some environments, and at certain times. This calls for building adaptive environments that can be personalized for the user and task to provide appropriate XAI methods, cognitive forcing functions, and other AI advice interventions. Such environments could improve AI-assisted decision making and help elicit more diverse data for training robust NLP models. However, in order to develop these environments, we need a deeper understanding of how environmental factors interact with user characteristics and different XAI methods. Future work should explore other pressures such as competition and multi-tasking, alongside XAI and measures of user personality and behavior, e.g., how they cope with stress. Additional work investigating the influence of performance pressure on tasks where humans have higher than random chance performance should also be investigated, as the literature suggests performance pressure can improve performance, despite our results suggesting minimal influence.

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

7 Conclusion

AI assistance is already used in the real world for tasks with high and low risk, and designing adaptive AI assistants that are domain-specific requires understanding the different factors influencing how humans use AI advice. We investigated how pressure influences the use of AI advice. Using deceptive review classification, crowdworkers with little expertise or personal motivation requirements, and two different XAI techniques, we observed complex effects on how pressure influences AI advice usage. Pressure and XAI interactions could both improve and decrease appropriate AI advice usage along the two dimensions of appropriate AI reliance and appropriate self-reliance. While performance pressure had minor effects, time pressure had a strong negative effect. Our results contribute to the body of literature investigating how pressure influences AI-assisted decision making. We note the relevance of these findings in AI annotation projects in particular; our work motivates continued research on the effects of pressure on AI assistance in varied environments while taking into account individual differences.

815

816

817

818

819

820

821

822

823

824

825

826

773

774

Limitations

724

726

727

728

730

731

732

739

740

741

742

743

761

762

763

767

770

771

772

Our experiments in manipulating environmental pressure used an inherently low-stakes, challenging task completed by crowdworkers. As a relatively young area of inquiry, it would be important to investigate how pressure and XAI influence changes between laypeople and experts. We expect different behavior because expertise can influence information seeking behavior (Cathy C. Durham and McLeod, 2000) and buffer stress responses (Matthews et al., 2019).

The literature has yet to find consensus on how people behave in AI-assisted decision making. Our LIME-baseline results contrast with those of Schemmer et al. (2023), despite using a similar interface and pool of crowdworkers. Understanding AI-assisted decision making is a complex and challenging endeavor, and the rapid adoption of AI assistance further motivates research in this area.

Ethical Considerations

Our work investigated how pressure and explain-744 able AI influence AI-assisted decision making, and 745 our results show how trust in the AI system plays 746 an important role in agreeing with AI advice. A 747 malicious actor could design a system to increase AI trust in order to persuade others to agree with AI 749 advice, against their better judgment, and overriding beneficial influences of pressure that persuade 751 people to be more skeptical and careful. For example, 10 of our participants stated they found the AI 753 advice agreed with their judgments for the first few 754 instances, leading them to rely on the AI advice more at later stages. Imposing a time limit on the decision could further convince people to rely on the AI advice because relying on the AI is an easy 758 method to cope with time stress.

References

- Robert H Ashton. 1990. Pressure and performance in accounting decision settings: Paradoxical effects of incentives, feedback, and justification. *Journal of Accounting Research*, 28:148–180.
- Yannick A Balk, Marieke A Adriaanse, Denise TD De Ridder, and Catharine Evers. 2013. Coping under pressure: Employing emotion regulation strategies to enhance performance under pressure. *Journal of Sport and Exercise Psychology*, 35(4):408–418.
- Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole

exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings* of the 2021 CHI Conference on Human Factors in Computing Systems, CHI '21, New York, NY, USA. Association for Computing Machinery.

- Daniel Ben David, Yehezkel S Resheff, and Talia Tron. 2021. Explainable ai and adoption of financial algorithmic advisors: an experimental study. In *Proceedings of the 2021 AAAI/ACM Conference on AI*, *Ethics, and Society*, pages 390–400.
- Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on ai in aiassisted decision-making. *Proceedings of the ACM on Human-computer Interaction*, 5(CSCW1):1–21.
- Adrian Bussone, Simone Stumpf, and Dympna O'Sullivan. 2015. The role of explanations on trust and reliance in clinical decision support systems. In 2015 International Conference on Healthcare Informatics, pages 160–169.
- Madalina Busuioc. 2021. Accountable artificial intelligence: Holding algorithms to account. *Public administration review*, 81(5):825–836.
- Noah Castelo, Maarten W. Bos, and Donald R. Lehmann. 2019. Task-dependent algorithm aversion. *Journal of Marketing Research*, 56(5):809–825.
- June M. L. Poon Cathy C. Durham, Edwin A. Locke and Poppy L. McLeod. 2000. Effects of group goals and time pressure on group efficacy, informationseeking strategy, and performance. *Human Performance*, 13(2):115–138.
- Karan R. Chhabra, Kathryn I. Pollak, Stephanie J. Lee, Anthony L. Back, Roberta E. Goldman, and James A. Tulsky. 2013. Physician communication styles in initial consultations for hematological cancer. *Patient Education and Counseling*, 93(3):573–578.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7282–7296, Online. Association for Computational Linguistics.
- Lawrence A. Crosby, Kenneth R. Evans, and Deborah Cowles. 1990. Relationship quality in services selling: An interpersonal influence perspective. *Journal of Marketing*, 54(3):68–81.
- Steven E Dilsizian and Eliot L Siegel. 2014. Artificial intelligence in medicine and cardiac imaging: harnessing big data and advanced computing to provide personalized medical diagnosis and treatment. *Current cardiology reports*, 16:1–8.

830

- 831

- 847
- 850 851 852
- 854
- 855

864

- 870

874

875 876 877

878 879

- Patricia M. Doney and Joseph P. Cannon. 1997. An examination of the nature of trust in buyer-seller relationships. Journal of Marketing, 61(2):35-51.
- Loïc Duron, Alexis Ducarouge, André Gillibert, Julia Lainé, Christian Allouche, Nicolas Cherel, Zekun Zhang, Nicolas Nitche, Elise Lacave, Aloïs Pourchot, et al. 2021. Assessment of an ai aid in detection of adult appendicular skeletal fractures by emergency physicians and radiologists: a multicenter crosssectional diagnostic study. Radiology, 300(1):120-129.
- Michael W Eysenck. 2012. Anxiety and cognitive performance. HANDBOOK OF PSYCHOLOGY OF EMOTIONS, page 87.
- Jürgen Fleiß, Elisabeth Bäck, and Stefan Thalmann. 2024. Mitigating algorithm aversion in recruiting: A study on explainable ai for conversational agents. SIGMIS Database, 55(1):56-87.
- Shankar Ganesan. 1994. Determinants of long-term orientation in buyer-seller relationships. Journal of Marketing, 58(2):1-19.
- Lior Gazit, Ofer Arazy, and Uri Hertz. 2023. Choosing between human and algorithmic advisors: The role of responsibility sharing. Computers in Human Behavior: Artificial Humans, 1(2):100009.
- David Gefen, Elena Karahanna, and Detmar W. Straub. 2003. Trust and tam in online shopping: an integrated model. MIS Q., 27(1):51-90.
- Ben Green and Yiling Chen. 2019. The principles and limits of algorithm-in-the-loop decision making. Proc. ACM Hum.-Comput. Interact., 3(CSCW).
- Nina Grgić-Hlača, Claude Castelluccia, and Krishna P Gummadi. 2022. Taking advice from (dis) similar machines: the impact of human-machine similarity on machine-assisted decision-making. In Proceedings of the AAAI Conference on Human Computation and Crowdsourcing, volume 10, pages 74-88.
- Nikita Haduong, Alice Gao, and Noah A. Smith. 2023. Risks and NLP design: A case study on procedural document QA. In Findings of the Association for Computational Linguistics: ACL 2023, pages 1248-1269, Toronto, Canada. Association for Computational Linguistics.
- Robin Hale, Jeremy J. Piggott, and Stephen E. Swearer. 2017. Describing and understanding behavioral responses to multiple stressors and multiple stimuli. Ecology and Evolution, 7(1):38–47.
- Lew Hardy, Richard Mullen, and Nikki Martin. 2001. Effect of task-relevant cues and state anxiety on motor performance. Perceptual and motor skills, 92(3):943-946.
- Jennifer Henderson, Maria Kavussanu, Andrew Cooke, and Christopher Ring. 2024. Some pressures are more equal than others: Effects of isolated pressure

on performance. Psychology of Sport and Exercise, 72:102592.

881

882

883

884

885

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

- Charles A. Holt and Susan K. Laury. 2002. Risk aversion and incentive effects. The American Economic Review, 92(5):1644-1655.
- Maia Jacobs, Melanie F. Pradier, Thomas H. McCoy, Roy H. Perlis, Finale Doshi-Velez, and Krzysztof Z. Gajos. 2021. How machine-learning recommendations influence clinician treatment selections: the example of the antidepressant selection. Translational psychiatry, 11(1).
- Jinglu Jiang, Surinder Kahai, and Ming Yang. 2022. Who needs explanation and when? juggling explainable ai and user epistemic uncertainty. International Journal of Human-Computer Studies, 165:102839.
- Deborah G Johnson. 2021. Algorithmic accountability in the making. Social Philosophy and Policy, 38(2):111-127.
- S Mo Jones-Jang and Yong Jin Park. 2023. How do people react to ai failure? automation bias, algorithmic aversion, and perceived controllability. Journal of Computer-Mediated Communication, 28(1):zmac029.
- Edgar E. Kausel, Satoris S. Culbertson, Pedro I. Leiva, Jerel E. Slaughter, and Alexander T. Jackson. 2015. Too arrogant for their own good? why and when narcissists dismiss advice. Organizational Behavior and Human Decision Processes, 131:33-50.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. Dynabench: Rethinking benchmarking in NLP. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4110-4124, Online. Association for Computational Linguistics.
- Kundan Krishna, Prakhar Gupta, Sanjana Ramprasad, Byron Wallace, Jeffrey Bigham, and Zachary Lipton. 2023. USB: A unified summarization benchmark across tasks and domains. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 8826–8845, Singapore. Association for Computational Linguistics.
- Kundan Krishna, Sanjana Ramprasad, Prakhar Gupta, Byron C Wallace, Zachary C Lipton, and Jeffrey P Bigham. 2024. Genaudit: Fixing factual errors in language model outputs with evidence. arXiv preprint arXiv:2402.12566.
- Kathryn Ann Lambe, Gary O'Reilly, Brendan D. Kelly, and Sarah Curristan. 2016. Dual-process cognitive interventions to enhance diagnostic reasoning: a systematic review. BMJ Qual Saf, 25:808-820.

Benjamin Charles Germain Lee, Doug Downey, Kyle Lo, and Daniel S Weld. 2023. Limeade: From ai explanations to advice taking. *ACM Transactions on Interactive Intelligent Systems*, 13(4):1–29.

937

938

941

943

947

951

952

954

959

963

964

965

969

970

972

974

975

976

977

978

981

991

- Bo-Ru Lu, Nikita Haduong, Chia-Hsuan Lee, Zeqiu Wu, Hao Cheng, Paul Koester, Jean Utke, Tao Yu, Noah A Smith, and Mari Ostendorf. 2024. Does collaborative human–Im dialogue generation help information extraction from human–human dialogues? In *First Conference on Language Modeling*.
 - Gerald Matthews, Ryan W. Wohleber, and Jinchao Lin. 2019. 490stress, skilled performance, and expertise: Overload and beyond. In *The Oxford Handbook of Expertise*. Oxford University Press.
 - Carolyn McKay. 2020. Predicting risk in criminal procedure: actuarial tools, algorithms, ai and judicial decision-making. *Current Issues in Criminal Justice*, 32(1):22–39.
 - Kimberly K. Merriman and John R. Deckop. 2007. Loss aversion and variable pay: a motivational perspective. *The International Journal of Human Resource Management*, 18(6):1026–1041.
 - Katelyn Morrison, Philipp Spitzer, Violet Turri, Michelle Feng, Niklas Kühl, and Adam Perer. 2024.
 The impact of imperfect xai on human-ai decisionmaking. *Proceedings of the ACM on Human-Computer Interaction*, 8(CSCW1):1–39.
 - Richard Mullen, Lew Hardy, and Andrew Tattersall. 2005. The effects of anxiety on motor performance: A test of the conscious processing hypothesis. *Journal of Sport and Exercise Psychology*, 27(2):212– 225.
 - Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4885–4901, Online. Association for Computational Linguistics.
 - Arne Nieuwenhuys and Raôul RD Oudejans. 2012. Anxiety and perceptual-motor performance: toward an integrated model of concepts, mechanisms, and processes. *Psychological research*, 76:747–759.
 - Charles A. O'Reilly and Nicholas Hall. 2021. Grandiose narcissists and decision making: Impulsive, overconfident, and skeptical of experts-but seldom in doubt. *Personality and Individual Differences*, 168:110280.
 - Myle Ott, Claire Cardie, and Jeffrey T. Hancock. 2013.
 Negative deceptive opinion spam. In Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 497–501, Atlanta, Georgia. Association for Computational Linguistics.

Myle Ott, Yejin Choi, Claire Cardie, and Jeffrey T. Hancock. 2011. Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 309–319, Portland, Oregon, USA. Association for Computational Linguistics. 992

993

994

995

996

997

998

999

1000

1001

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1033

1034

1035

1036

1038

1039

1040

1041

1042

- Cecilia Panigutti, Andrea Beretta, Fosca Giannotti, and Dino Pedreschi. 2022. Understanding the impact of explanations on advice-taking: a user study for ai-based clinical decision support systems. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–9.
- Charvi Rastogi, Yunfeng Zhang, Dennis Wei, Kush R Varshney, Amit Dhurandhar, and Richard Tomsett. 2022. Deciding fast and slow: The role of cognitive biases in ai-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW1):1–22.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings* of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, page 1135–1144, New York, NY, USA. Association for Computing Machinery.
- David Ronayne and Daniel Sgroi. 2019. Ignoring good advice. University of Warwick, Warwick Economics Research Papers Series, 1150.
- Max Schemmer, Niklas Kuehl, Carina Benz, Andrea Bartos, and Gerhard Satzger. 2023. Appropriate reliance on ai advice: Conceptualization and the effect of explanations. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, IUI '23, page 410–422, New York, NY, USA. Association for Computing Machinery.
- Helen Smith. 2021. Clinical ai: opacity, accountability, responsibility and liability. *Ai & Society*, 36(2):535–545.
- Janet A Sniezek and Timothy Buckley. 1995. Cueing and cognitive conflict in judge-advisor decision making. *Organizational behavior and human decision processes*, 62(2):159–174.
- Mark Steyvers and Aakriti Kumar. 2023. Three challenges for ai-assisted decision-making. *Perspectives on Psychological Science*, page 17456916231181102.
- Mike Stoker, Ian Maynard, Joanne Butt, Kate Hays, and Paul Hughes. 2019. The effect of manipulating individual consequences and training demands on experiences of pressure with elite disability shooters. *The sport psychologist*, 33(3):221–227.
- Siddharth Swaroop, Zana Buçinca, Krzysztof Z. Gajos,
and Finale Doshi-Velez. 2024. Accuracy-time trade-
offs in ai-assisted decision making under time pres-
sure. In Proceedings of the 29th International Con-
ference on Intelligent User Interfaces, IUI '24, page1044
10451046
10471046
1046

- disregards incorrect AI advice. $RAIR = \frac{\sum_{i=0}^{N} CSR_i}{\sum_{i=0}^{N} IA_i}$ ity Questionnaire–276 (MPQ-276) test booklet. Uni-(2)1096 1. CSR: 1 if the initial judgment agrees is cor-Somin Wadhwa, Silvio Amir, and Byron C Wallace. 1097 2024. Investigating mysteries of CoT-augmented disrect, the AI advice is incorrect, and the final tillation. In Proceedings of the 2024 Conference on judgment is correct; 0 otherwise 1099 Empirical Methods in Natural Language Processing, pages 6071-6086, Miami, Florida, USA. Association 2. IA: 1 if the initial judgment is correct and the 1100 AI advice is incorrect; 0 otherwise 1101 Xiuxin Wang and Xiufang Du. 2018. Why does advice A.2 Presurvey 1102 discounting occur? the combined roles of confidence The presurvey contained demographics questions 1103 and the 10-item IPIP adaptation of the Tellegen 1104 Ilan Yaniv and Eli Kleinberger. 2000. Advice taking in decision making: Egocentric discounting and reputa-(1995/2003) Multiple Personality Questionnaire 1105 tion formation. Organizational behavior and human (MPQ). Questions: 1106 1. Indicate your age range 1107 2. Indicate your race 1108 3. Indicate your gender 1109 4. What is your highest level of education com-1110 pleted? 1111 5. What is your native language? 1112 6. Describe your proficiency in other languages 1113 7. How familiar are you with the task of deciding 1114 whether a review is genuine or deceptive? (5-1115 point Likert) 1116 8. Rate your level of agreement with the follow-1117 ing statements (5-point Likert): (1) I trust ar-1118 tificial intelligence (AI); (2) I use AI to help 1119 me with my work 1120 9. MPQ 1121 MPQ consists of 10 questions answered on a 4-1122 (1)point Likert scale. Subjects are asked to rate their 1123 level of agreement with each statement. 1124 · I would never go hang gliding or bungee jump-1125 ing. 1126 • I would never make a high-risk investment. 1127 I avoid dangerous situations. 1128 • I seek danger. 1129
 - I am willing to try anything once. 1130
 - 12
- 1. CAIR: 1 if the initial judgment disagrees with the ground truth, the AI advice is correct, and the final judgment agrees with the ground truth: 0 otherwise
 - 2. CA: 1 if the initial judgment is incorrect and the AI advice is correct, regardless of the final judgment; 0 otherwise

where

XAI condition were generated using ChatGPT with default settings on Jan 20, 2025, with the prefix prompt: In less than 50 words, explain why the following review is a [genuineldeceptive] review. A.1 RAIR and RSR

Relative AI reliance (RAIR) and relative selfreliance (RSR) are defined by Schemmer et al. (2023), which we reproduce below for reference. Subjects complete a prediction task with N instances $x_i \in X$ with ground truth labels $y_i \in Y$. RAIR is the ratio of cases where the human correctly changes their decision to follow AI advice.

 $RAIR = \frac{\sum_{i=0}^{N} CAIR_i}{\sum_{i=0}^{N} CA_i}$

ing: methodologies and literature review. Journal of *Multi-Criteria Decision Analysis*, 11(4-5):167–186. Methods Α The natural language explanations for the GenAI

decision processes, 83(2):260–281. Constantin Zopounidis and Michael Doumpos. 2002. Multi-criteria decision aid in financial decision mak-

and trust. Frontiers in psychology, 9:2381.

138-154, New York, NY, USA. Association for Com-

puting Machinery.

versity of Minnesota Press.

for Computational Linguistics.

1049

1050

1054 1055

1057

1059

1063

1066

1069

1071

1072

1073

1074

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1093

Auke Tellegen. 1995/2003. Multidimensional Personal-

RSR is the ratio of cases where the human correctly



(a) The initial decision

(b) Second page of the judgment task displaying AI advice alongside LIME explanation

Figure 3: The user interface in payment-LIME condition. The user first encounters the review in (a) and makes a judgment. Then they receive AI advice (b) in the form of a model prediction and need to make their judgment again.

- I do dangerous things.
- I enjoy being reckless
 - I seek adventure
 - I take risks

1133

1134

1135

1136

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1152

• I do crazy things

A.3 Postsurvey

- 1137 Postsurvey questions:
 - Trust in AI Advisor: Rate your level of agreement with the following statements on a 7-point Likert scale.
 I think I can trust the AI Advisor; 2) The AI advisor can be trusted to provide reliable support; 3) I trust the AI advisor to keep my best interests in mind; 4) In my opinion, the AI advisor is trustworthy
 - 2. HL survey
 - 3. Describe how you determined whether a review was genuine or deceptive.

1148The Trust in AI Advisor questions were sourced1149from Schemmer et al. (2023), who sourced them1150from Crosby et al., 1990; Doney and Cannon, 1997;1151Ganesan, 1994 and Gefen et al., 2003.

A.3.1 Holt-Laury Survey

1153The Holt-Laury survey measures how risk averse1154a respondent is by asking them to make a decision1155between pairs of gambles with a "safe" and "risky"1156choice. For example, in the first question, the par-1157ticipant can select the safe choice of having 1/101158chance to earn \$2 and 9/10 chance of earning \$1.60,

or the risky choice of having 1/10 chance to earn1159\$3.85 and 9/10 chance of earning \$0.10. One of1160the choices is randomly selected and the gamble1161played to determine the bonus to the respondent.1162

Safe choice	Risky choice
1/10 of \$2.00, 9/10 of \$1.60	1/10 of \$3.85, 9/10 of \$0.10
2/10 of \$2.00, 8/10 of \$1.60	2/10 of \$3.85, 8/10 of \$0.10
3/10 of \$2.00, 7/10 of \$1.60	3/10 of \$3.85, 7/10 of \$0.10
4/10 of \$2.00, 6/10 of \$1.60	4/10 of \$3.85, 6/10 of \$0.10
5/10 of \$2.00, 5/10 of \$1.60	5/10 of \$3.85, 5/10 of \$0.10
6/10 of \$2.00, 4/10 of \$1.60	6/10 of \$3.85, 4/10 of \$0.10
7/10 of \$2.00, 2/10 of \$1.60	7/10 of \$3.85, 2/10 of \$0.10
8/10 of \$2.00, 3/10 of \$1.60	7/10 of \$3.85, 3/10 of \$0.10
8/10 of \$2.00, 2/10 of \$1.60	8/10 of \$3.85, 2/10 of \$0.10
9/10 of \$2.00, 1/10 of \$1.60	9/10 of \$3.85, 1/10 of \$0.10
10/10 of \$2.00, 0/10 of \$1.60	10/10 of \$3.85, 0/10 of \$0.10

Table 1: Holt-Laury survey

B Regression Tables

The random effects model regression tables are1164detailed below.1165

DV: RAIR	Coef.	Std.Err.	Z	P> z	[0.025	0.975]
Intercept	0.548	0.059	9.310	0.000	0.432	0.663
Pressure	-0.022	0.012	-1.761	0.078	-0.046	0.002
Explanation Condition	-0.008	0.024	-0.334	0.738	-0.056	0.040
MPQ	-0.075	0.041	-1.812	0.070	-0.156	0.006
Holt-Laury	-0.059	0.043	-1.367	0.172	-0.143	0.025
Advisor Trust	-0.193	0.061	-3.158	0.002**	-0.312	-0.073
Trust in AI	-0.016	0.018	-0.923	0.356	-0.051	0.018
Age	-0.023	0.020	-1.141	0.254	-0.062	0.016
Race	0.003	0.009	0.380	0.704	-0.014	0.020
Gender	0.040	0.034	1.177	0.239	-0.027	0.107
Education	0.014	0.007	1.949	0.051	-0.000	0.029
Subject Group	0.083	0.038				

Table 2: DV: RAIR. * p < 0.1, ** p < 0.05, ***p < 0.01

DV: RSR	Coef.	Std.Err.	Z	P> z	[0.025	0.975]
Intercept	0.618	0.062	10.017	0.000	0.497	0.739
Pressure	-0.043	0.012	-3.595	0.000	-0.066	-0.020
Explanation Condition	-0.069	0.026	-2.686	0.007**	-0.119	-0.019
MPQ	0.051	0.044	1.161	0.246	-0.035	0.136
Holt-Laury	0.044	0.045	0.974	0.330	-0.045	0.133
Advisor Trust	0.164	0.064	2.559	0.011*	0.038	0.290
Trust in AI	0.005	0.018	0.294	0.768	-0.031	0.042
Age	0.014	0.021	0.666	0.506	-0.027	0.055
Race	-0.004	0.009	-0.471	0.638	-0.022	0.013
Gender	-0.035	0.036	-0.986	0.324	-0.106	0.035
Education	-0.005	0.008	-0.589	0.556	-0.020	0.011
Subject Group	0.098	0.044				

Table 3: DV: RSR. * p < 0.1, ** p < 0.05, ***p < 0.01

DV: Overall Accuracy	Coef.	Std.Err.	Z	P> z	[0.025	0.975]
Intercept	0.598	0.018	33.481	0.000	0.563	0.632
Pressure	-0.016	0.006	-2.721	0.007**	-0.028	-0.005
Explanation Condition	-0.036	0.007	-5.096	0.000	-0.050	-0.022
MPQ	0.000	0.012	0.040	0.968	-0.023	0.024
Holt-Laury	-0.006	0.013	-0.489	0.625	-0.031	0.018
Advisor Trust	0.005	0.018	0.260	0.795	-0.030	0.040
Trust in AI	-0.002	0.005	-0.383	0.701	-0.012	0.008
Age	-0.006	0.006	-1.036	0.300	-0.018	0.005
Race	-0.002	0.003	-0.774	0.439	-0.007	0.003
Gender	-0.006	0.010	-0.622	0.534	-0.026	0.013
Education	0.005	0.002	2.362	0.018*	0.001	0.009
Subject Group	0.003	0.007				

Table 4: DV: Accuracy. * p < 0.1, ** p < 0.05, ***
 p < 0.01

1166	C Licenses
1167	The Deceptive Opinion Spam Corpus v1.4 (Ott
1168	et al., 2011, 2013) is licensed under Creative Com-
1169	mons Attribution-NonCommercial-ShareAlike 3.0
1170	Unported License. The outputs from ChatGPT, a
1171	large language model from OpenAI, are copyright
1172	free.
1173	D AI Use
1174	ChatGPT was used to help with data transforma-

tion.