

# The Effect of Feedback on News-Verification Demand: Experimental Evidence\*

[\(Check the latest version here.\)](#)

By DARIO TRUJANO-OCHOA<sup>†</sup> AND JOSE GLORIA<sup>‡</sup>

*This study is the first to elicit the value individuals place on verifying real headlines, focusing on the role of feedback in shaping verification behavior and confidence. In a lab experiment, 184 participants classified headlines as accurate or fake, reporting their willingness to pay (WTP) for a perfect signal under three feedback conditions: control, individual, and group. Results show that explicit feedback on others' accuracy reduces the perceived value of verification, potentially discouraging fact-checking. The findings highlight how feedback influences verification demand. By addressing gaps in existing literature and focusing on Mexico, this study provides novel insights into the design of interventions to combat misinformation.*

*JEL: C93, D83, D91*

*Keywords: Misinformation, Feedback, Verification, Willingness to Pay*

When an individual encounters a headline, they face a decision problem: should they invest time and effort to verify its accuracy before forming an opinion or acting on it? This decision often hinges on the perceived cost of verification versus the potential consequences of relying on misinformation. If individuals are informed about how difficult it is to identify false claims, they may reassess the value of verification. Knowing the challenge may increase their willingness to verify, as they recognize the higher risk of error, or conversely, it could discourage them if they perceive the task as too daunting, especially if the headline aligns with their pre-existing beliefs.

The spread of misinformation through digital platforms is faster and more profound than factual information (Vosoughi et al., 2018), and it is primarily people (not bots) who share misinformation inadvertently (Arin et al., 2023). Focusing on interventions at the individual decision-making level is essential to prevent

\* The authors want to thank the professors who discussed the early versions of the present project: Gary Charness, Cesi Cruz, Daniel Martin, Ignacio Esponda, Ryan Oprea, Sevgi Yuksel, Erik Eyster, and to the specialists and professors in Mexico who increased the discussion and understanding of the relevant problem of misinformation: Grisel Salazar, Daniel Moreno, Horacio Larreguy, Antonio Arechar, Pablo Soto, and Arturo Bouzas. We also thank AlianzaMX for the fellowship that allowed the authors to travel to Mexico to develop this research.

<sup>†</sup> UCSB, USA, [dariotrujanoochoa@ucsb.edu](mailto:dariotrujanoochoa@ucsb.edu).

<sup>‡</sup> UCLA, USA, [josegloria@ucla.edu](mailto:josegloria@ucla.edu).

people from making decisions based on false assumptions. Misinformation undermines the public’s ability to make informed decisions, erodes trust in institutions, and fuels polarization, threatening democratic societies.

Most existing research focuses on the effects of exogenous interventions and overlooks the intrinsic value individuals place on verifying the accuracy of the information they receive. This value ultimately determines the time and effort they are willing to invest in checking the veracity of headlines. In this paper, we refer to verification as searching for further information to reveal if a statement is true or false.<sup>1</sup> Individual verification is essential to fight misinformation, increasing the effectiveness of measures that point out false statements like debunking and fact-checking.<sup>2</sup> Thus, verification complements other efforts to fight misinformation.<sup>3</sup>

Overconfidence has been identified as a factor contributing to the spread of misinformation, as overconfident individuals are less likely to invest time and effort in verifying the accuracy of information. However, evidence supporting this relationship primarily comes from survey data. In this paper, we formalize this reasoning with an explicit framework for the decision problem of information classification, where individuals highly confident in their ability to identify false information are predicted to reduce their demand for verification.

The primary objective was to assess how a feedback intervention impacts individuals’ value on verification. While interventions like debunking, inoculation, and digital literacy provide indirect feedback on accuracy when classifying headlines, this research directly measures the demand for verification in a lab experiment, examining how explicit feedback alters this valuation. Some evidence suggests that feedback reduces overconfidence, which motivated us to explore its potential effects in this study. Still, the literature is far from a consensus on the impact of feedback. This is the first study of the effect of feedback on one’s confidence in spotting false information and willingness to pay to verify it.

We designed a classification-verification game implemented in a lab experiment to examine how feedback influences individuals’ demand for verifying news headlines. The experiment involved 184 participants from Mexican universities, who were tasked with classifying 50 headlines as accurate or fake across five blocks. Participants reported their willingness to pay (WTP) for a “perfect signal” that verified a headline’s accuracy. Feedback was varied among three groups: no feedback, individual feedback, and group feedback. The study employed a

<sup>1</sup>Not every statement and information can be classified this way. Normative statements (e.g., “Headlines should prioritize impactful language to capture readers’ attention.”) are examples. This paper focuses on headlines easily classified as accurate or false.

<sup>2</sup>In this paper, we distinguish fact-checking from verification by defining fact-checking as a policy aimed at identifying and disseminating information about fake news, whereas verification involves an individual’s proactive choice to seek additional details.

<sup>3</sup>For example, Facebook (Meta) collaborates with fact-checking organizations to tag false information (<https://transparency.meta.com/policies/community-standards/misinformation>), and X uses CommunityNotes (<https://x.com/communitynotes?lang=en>) to let the users decide through an algorithm over their ratings.

decision-making framework to link WTP to the perceived value of verification, using controlled conditions to isolate the effects of feedback on participants’ accuracy, confidence, and demand for verification. The experiment also incorporated incentivized tasks to ensure meaningful engagement, and it measured the impact of headlines’ political content and participants’ political leaning on participants’ verification behaviors.

To our knowledge, this is the first study that elicits the value of verifying actual headlines. It contributes to the literature by examining the impact of feedback on both confidence levels and the perceived value of verification, providing evidence that feedback on the accuracy of other participants can reduce this value. By focusing on headlines from Mexico—the largest Spanish-speaking country where misinformation and polarization have recently increased—this research addresses a critical gap in the literature, which has predominantly concentrated on interventions studied in the U.S. and Europe (Kozyreva et al., 2024; Bateman and Jackson, 2024). The study’s experimental design controls for participant heterogeneity, ensuring that its results are broadly applicable to diverse populations. Moving beyond the typical focus on exogenous interventions, this paper innovatively explores how feedback mechanisms influence individuals’ active demand for verification, offering insights into how such mechanisms can enhance the effectiveness of fact-checking in combating misinformation.

## I. Previous Literature

People’s overconfidence has been shown to impact their ability to discern fake news, leading to greater engagement with false information. According to Lyons et al. (2021), individuals tend to overestimate their ability to distinguish between real and fake news. Overconfident individuals are more likely to believe in and share fake news and are unaware of their limitations in identifying misinformation. This overconfidence exacerbates the spread of false information, especially when the misinformation aligns with their political or ideological beliefs. Similarly, Pennycook and Rand (2020) found that overconfidence in one’s cognitive abilities correlates with a higher likelihood of accepting false claims as accurate. Ortoleva and Snowberg (2015) found that overconfidence due to correlation neglect leads to higher polarization. Together, these studies highlight the role of overconfidence in the dissemination of fake news by decreasing verification and increasing misinformation sharing.

Encouraging fact-checking and verification behaviors has been effective in combating misinformation. Reviews by Kozyreva et al. (2024) and Bateman and Jackson (2024) noted that promoting media literacy, fact-checking, and labeling content are vital tools to counter misinformation. However, Aslett et al. (2024) found that verification tools, like search engines, could backfire by amplifying misinformation when individuals are overconfident in their search skills. This research indicates that while verification behaviors are crucial, the limitations of tools like search engines must be considered in curbing misinformation effectively.

Other interventions, such as accuracy prompts, encourage users to pause and assess content accuracy before sharing. For instance, Pennycook et al. (2021) introduced an “accuracy nudge” intervention that reduced fake news spread by prompting users to reflect on accuracy, highlighting the positive impact of slight cognitive adjustments on online behavior. Similarly, Pennycook and Rand (2022) focused on interventions promoting cognitive engagement to reduce misinformation dissemination, demonstrating that simple adjustments in thought processes significantly improve the quality of shared content.

The literature shows mixed results on feedback’s effects on overconfidence. While some incentivized studies have found feedback reduces overconfidence (Ferraro, 2005; Eberlein et al., 2011; Kogelnik, 2022), others in educational contexts report no effects (Pulford and Colman, 1997; Erat et al., 2022). Moreover, some studies reveal feedback’s asymmetric effects on motivated beliefs. Oprea and Yuksel (2022) observed that feedback increases subjective probabilities of outperforming others, while Thaler (2024) noted that individuals update beliefs asymmetrically when feedback reinforces ego-related worldviews. Kartal and Tyran (2022) showed that overconfident participants continue voting despite low-accuracy signals, though they did not directly measure demand for information. Additionally, Moore and Healy (2008) distinguished between overestimation, overplacement, and overprecision, noting that feedback often has minimal effects on recalibrating overconfidence, particularly for difficult tasks where individuals overestimate their performance.

This research aims to establish whether feedback impacts the demand for verification. Unlike prior studies focusing on the accuracy of misinformation distinction or verification strategies, this study offers experimental evidence on feedback’s causal effect on verification demand, a novel contribution to the misinformation literature.

## II. Hypotheses on the Effects of Feedback

This study is structured around three main hypotheses that explore the impact of feedback on willingness to pay (WTP) to verify the accuracy of the headlines and the confidence. Two other secondary hypotheses are included regarding the effects on accuracy. Under the assumption of rationality, participants’ WTP should be proportional to the value they give to information. These hypotheses are integrated into the methodology to test their validity in a controlled environment.<sup>4</sup>

**HYPOTHESIS 1:** *Participants are generally overconfident and will expect better performance in classifying headlines than what is reflected in the feedback they receive.*

<sup>4</sup>This study was preregistered in OSF: <https://osf.io/jxr82>. We changed the order of the hypothesis for exposition purposes, leaving the hypothesis about accuracy at the end.

Following the literature on overconfidence, this hypothesis suggests that participants tend to overestimate their classification accuracy before receiving feedback. In this experiment, participants' predictions about their classification accuracy are expected to be higher than the accuracy indicated by their actual performance, as revealed by the feedback.

**HYPOTHESIS 2:** *Participants' willingness to pay (WTP) for verification will be higher when they receive feedback on group classification accuracy than when they receive personal performance feedback.*

According to previous results on the asymmetric effect of feedback, participants who receive feedback on the performance of others could perceive this feedback as more informative and objective. As a result, they will place a larger value on the verification process and be willing to pay more to ensure the accuracy of the headlines they classify. The expectation is that group feedback, being less affected by motivated reasoning, will lead to a higher WTP as participants seek to mitigate the perceived difficulty of the task.

**HYPOTHESIS 3:** *Participants' willingness to pay (WTP) for verification is influenced by the political content of the headlines, with differing effects depending on whether the headline favors or opposes the current government.*

- 1) *Supporters of the Government: Lower WTP for verification of favorable headlines; higher WTP for verification of unfavorable headlines.*
- 2) *Opponents of the Government: Higher WTP for verification of favorable headlines; lower WTP for verification of unfavorable headlines.*

When participants are presented with headlines that contain political content, it is hypothesized that their WTP for verification will vary depending on whether the headline aligns with their political beliefs. Specifically, if a headline is favorable to the current government, participants who support the government will likely have lower WTP for verification. This is because they are more inclined to accept information that aligns with their pre-existing beliefs without seeking further verification. Conversely, participants who oppose the current government may exhibit higher WTP for verification of favorable headlines, as they may be more skeptical of information that contradicts their beliefs and, thus, more motivated to confirm its accuracy.

On the other hand, for headlines unfavorable to the current government, supporters of the government may demonstrate higher WTP for verification, driven by a desire to challenge or disprove information that opposes their political views. Opponents of the government, however, may show lower WTP for verification of unfavorable headlines, as they may be more likely to accept information that aligns with their negative views of the government without the need for additional confirmation.

### A. Secondary Hypothesis

**HYPOTHESIS 4:** *Feedback on personal performance improves the accuracy of headline classification compared to no feedback.*

**HYPOTHESIS 5:** *Feedback on personal performance improves the accuracy of headline classification compared to no feedback.*

According to these hypotheses, giving participants feedback on their performance in classifying headlines will lead to higher accuracy in future classification tasks. The expectation is that personal feedback will encourage participants to adjust their behavior, leading to improved accuracy. This is related to previous studies showing that people improve when they are prompted to the importance of accuracy (Pennycook et al., 2021).

## III. Decision-Making in the Classification of Headline Accuracy

This section describes the agent’s decision-making problem of classifying, verifying, and reclassifying a headline as accurate (true) or fake (false). Classifying is a signal detection problem, and purchasing additional signals on this decision requires calculating the expected value of sample information. The instrumental value of information is assumed to be equal to the willingness to pay for verification. The preset setting provides a framework for understanding the decision-making process when deciding whether to verify headlines.

### A. Problem Setup Without Purchasing a Signal

Consider an agent tasked with classifying headlines as accurate ( $t$ ) or fake ( $f$ ) ( $c \in \{t, f\}$ ). The state of the world is  $\omega \in \Omega = \{T, F\}$ , with the prior probability of encountering a fake headline denoted by  $P(\omega = F) = p_f$ .<sup>5</sup> Consequently, the prior probability of encountering an accurate headline is  $1 - p_f$ .

The agent’s utility for correctly classifying a headline as accurate is  $U_T$ , and for correctly classifying a headline as fake is  $U_F$ ; . Conversely, the utility is  $U_{TF} < U_T$  for misclassifying a fake headline as accurate and  $U_{FT} < U_F$  for misclassifying an accurate headline as fake. The condition is case-insensitive for evaluating the correct classification (i.e.,  $c = \omega$  means a correct classification).

The probability of classifying correctly the headline is determined by  $P(c = t|T)$  and  $P(c = f|F)$  with  $1 < \frac{P(c=t|T)}{P(c=t|F)}$  and  $1 < \frac{P(c=f|F)}{P(c=f|T)}$  to assure that the initial classification  $c$  is informative in the sense that the initial classification  $c$  gives information relative to the prior probability of each state  $\omega$ .<sup>6</sup> This is stated

<sup>5</sup>We simplify  $P(\omega = F)$  to  $P(F)$ . For  $c, s \in \Omega$ , we specify.

<sup>6</sup>We are assuming here that the classification is a signal to the same agent without considering the content of a headline  $h$  which is most likely multidimensional. This classification process also follows an optimization process where  $c = \omega \iff \frac{P(\omega|h)}{P(\omega|\bar{\omega} \neq \omega|h)} > \frac{U_T - U_{FA}}{U_F - U_{TF}}$ . However, following the objectives of the present research, we focus on analyzing the informativeness of the initial classification  $c$  in proposition 1 without analyzing the properties of the headlines or the payoffs.

formally in the following proposition.

**PROPOSITION 1:** *Informativeness of the initial classification  $c$ .*

$$1 < \frac{P(c=\omega|\omega)}{P(c=\omega|\bar{\omega}\neq\omega)} \text{ if and only if } P(\omega|c \neq \omega) < P(\omega) < P(\omega|c = \omega)$$

Notice that commonly found assumptions  $0.5 < P(s = t|T) = q_t$  and  $0.5 < P(s = f|F) = q_f$  are sufficient to make the signal  $S$  informative according to proposition 1. The proof of this proposition can be found in the appendices.

The expected utilities when a headline is classified as accurate,  $EU_{\text{no signal}}(t)$ , and as fake,  $EU_{\text{no signal}}(f)$ , are given respectively by the equations:

$$EU_{\text{no signal}}(t) = P(T|t) \cdot U_T + P(F|t) \cdot U_{TF}$$

$$EU_{\text{no signal}}(f) = P(F|f) \cdot U_F + P(T|f) \cdot U_{FT}$$

The previous equations clearly show that if the headline's classification is informative, reading a headline is valuable because the expected utility is larger than considering the prior probabilities alone.

### B. Conditional WTP Analysis

This section showed the optimal willingness to pay (WTP) for signal  $S$  after the initial classification. The WTP is the maximum amount that an agent would pay to observe signal  $S$ .

The decision to purchase information happens after observing a headline once the agent has classified the signal. Therefore, the value of the signal depends on  $c$ . We are assuming that sequential information acquisition is optimal. The problem of sequential decision-making was stated in general by [Wald \(1947\)](#), and [Arrow et al. \(1949\)](#) analyzed how to learn from sequential information.

This section presents the condition that makes verifying the initial classification valuable. After initially classifying the headline as accurate ( $c = t$ ) or false ( $c = f$ ), the agent can reclassify the headline  $r \in \{t, f\}$  based on the signal's realization  $s \in \{t, f\}$ . Let's consider first a valuable signal  $S$  with the conditions in proposition 2.

**DEFINITION 1:** *A signal is valuable if  $EU(r = s) \geq EU(r = c)$*

The informativeness of the signal is determined by  $P(s = t|T) = q_t$ . We need a strong enough signal  $S$  so that the signal is valuable and the optimal decision to reclassify is to follow the signal ( $r = s \in \{t, f\}$ ). Also, we assume that the signal realization  $s$  is independent of the previous classification  $c$  conditional on the state of the world  $\omega \in \{T, F\}$  (i.e.  $P(s|\omega, c) = P(s|\omega)$ ).

**PROPOSITION 2:** *Conditions for Valuable Signal*

*A signal  $S$  is valuable if and only if*

$$\frac{P(\omega|s = \omega, c \neq \omega)}{P(\tilde{\omega}|s = \omega, c \neq \omega)} < \frac{U_\omega - U_{\tilde{\omega}\omega}}{U_{\tilde{\omega}} - U_{\omega\tilde{\omega}}} \equiv U_\omega$$

with  $\tilde{\omega} \in \Omega, \tilde{\omega} \neq \omega$ .

We are also assuming that the initial classification is valuable and therefore follows the analogous condition  $\frac{P(\omega|c=\omega)}{P(\tilde{\omega}|c=\omega)} < U_\omega$ .

By allowing the agent to update their classification based on the signal, we account for the dynamic decision-making process. The WTP to verify is derived by comparing the expected utility with the signal (considering reclassification) to the expected utility without the signal. This approach shows the impact of additional information on improving decision-making accuracy. The detailed mathematical steps and proofs are provided in the appendix. The expected utility of reclassification is calculated by updating the agent's posterior beliefs using Bayes' rule and comparing the expected utilities with and without reclassification.

For an agent tasked with classifying headlines as accurate or fake, a signal  $S$  indicating the state of the world must be sufficiently strong to ensure that the agent reclassifies based on this signal.

#### WTP EQUATION

Initially, the agent classifies a headline as either accurate ( $t$ ) or fake ( $f$ ). The agent updates their beliefs upon receiving a signal  $s$ , which can either confirm or contradict the initial classification. The posterior probabilities are calculated using Bayes' rule. For example, the posterior probability of the headline being accurate given the signal  $s = t$  and the initial classification  $c$  is:

$$P(T|s = t, c) = \frac{q_t \cdot P(T|c)}{q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)}$$

Similarly, the posterior probability of the headline being fake given the signal  $s = f$  and the initial classification  $c$  is:

$$P(F|s = f, c) = \frac{q_f \cdot P(F|c)}{(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)}$$

The agent's decision to reclassify based on the signal depends on the expected utilities. The expected utility of reclassification given the signal  $s = t$ , or  $s = f$ , are respectively:

$$EU_{\text{new classification}}(s = t, c) = P(T|s = t, c) \cdot U_T + P(F|s = t, c) \cdot U_{TF}$$

$$EU_{\text{new classification}}(s = f, c) = P(F|s = f, c) \cdot U_F + P(T|s = f, c) \cdot U_{FT}$$

The combined expected utility of updating the signal, considering both possible



signals, is:

$$\begin{aligned} EU_{\text{signal}}^{\text{update}}(c) &= P(s = t|c) \cdot EU_{\text{new classification}}(s = t, c) + \\ &\quad P(s = f|c) \cdot EU_{\text{new classification}}(s = f, c) \\ &= [q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)] \cdot [P(T|s = t, c) \cdot U_T + P(F|s = t, c) \cdot U_{TF}] + \\ &\quad [(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)] \cdot [P(F|s = f, c) \cdot U_F + P(T|s = f, c) \cdot U_{FT}] \end{aligned}$$

The WTP to verify the headline is calculated by comparing the expected utility with the signal to the expected utility without the signal:

$$V(c) = EU_{\text{signal}}^{\text{update}}(c) - EU_{\text{no signal}}(c)$$

### C. Simplifying Assumptions

Let's assume equal prior probabilities  $p_f = 0.5$  and equal utilities  $U_T = U_F = 1$  and  $U_{TF} = U_{FT} = 0$ . Also, assume that the prevalence of fake and accurate news is the same  $P(T) = P(F) = p_f = 0.5$ . These assumptions about the payoffs allow us to interpret the signal's value purely in probability terms related to its informativeness. Substituting these assumptions into the expected utility equations, we get:

$$EU_{\text{no signal}}(t) = P(T|t) \cdot 1 + P(F|t) \cdot 0 = P(T|t) = \frac{P(t|T)}{P(t|T) + P(t|F)}$$

$$EU_{\text{no signal}}(f) = P(F|f) \cdot 1 + P(T|f) \cdot 0 = P(F|f) = \frac{P(f|F)}{P(f|T) + P(f|F)}$$

And the expected utilities simplify to:

$$EU_{\text{new classification}}(s = t, c) = P(T|s = t, c)$$

$$EU_{\text{new classification}}(s = f, c) = P(F|s = f, c)$$

Finally, the combined expected utility of updating the signal, considering both possible signals, is:

$$\begin{aligned} EU_{\text{signal}}^{\text{update}}(c) &= [q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)] \cdot P(T|s = t, c) + \\ &\quad [(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)] \cdot P(F|s = f, c) \end{aligned}$$

### PERFECT SIGNAL

Here, we calculate the WTP considering the condition  $q_f = q_t = 1$ ; perfect signal. This assumption ensures that the signal is strong enough to follow even without the other simplifying assumptions and substantially simplifies the interpretation of  $EVSI(c)$ . For the case  $c = f$  and  $s = t$ :  $q_t \cdot P(T|f) > (1 - q_f) \cdot P(F|f) \iff$

$P(T|f) > 0$ . The case  $c = s$  and  $s = f$  requires  $P(F|t) > 0$ . Both conditions are satisfied by the construction of the problem. This assumption simplifies the expected utility of observing the signal  $S$ . Thus, the combined expected utility with the signal is:

$$EU_{\text{signal}}^{\text{update}}(c) = P(T|c) \cdot 1 + P(F|c) \cdot 1 = P(T|c) + P(F|c) = 1$$

Therefore,

$$(1) \quad V(c) = \begin{cases} 1 - P(T|t), c = t \\ 1 - P(F|f), c = f \end{cases}$$

The value of the signal  $S$  is the difference between the posterior probability of reclassifying correctly after observing the signal and the posterior probability of initially classifying correctly. Notice that if we change the payoff of a correct answer such that  $U_T = U_F > U_{TF} = U_{FT}$ , we only have to multiply the posterior probabilities difference by  $\pi = U_T - U_{TF}$  to get  $V(c)$ . Therefore, the willingness to pay to verify should be:

$$(2) \quad WTP(c) = \pi V(c)$$

The WTP decision is based solely on accuracy probability.

#### IV. Experimental Design: Classification-Verification Game

##### A. Overview

This experiment tests whether different types of feedback influence participants' accuracy in classifying information and their willingness to pay (WTP) for verification. Participants are tasked with categorizing headlines as accurate or fake and reporting their WTP for verification. They evaluated 50 headlines in five blocks, with feedback provided in one of three experimental conditions: control (no feedback), individual feedback, and group feedback. The design measures participants' classification accuracy, confidence, and willingness to pay for verification. All this is framed through the decision-making problem presented in section III with a perfect signal. Following the main task, participants complete a survey on demographics and political orientation.

##### B. Experimental Blocks and Feedback Treatments

Participants completed five blocks designed to measure classification accuracy, confidence, and demand for verification through WTP. In each block, participants received ten headlines in random order, which they were instructed to classify as either *true* ( $t$ ) or *false* ( $f$ ), with an equal prior probability ( $P(T) = P(F) = 0.5$ ) of each state. For each correctly classified headline, participants earned a utility

of 10 Mexican Pesos (MXN), regardless of whether the classification was *true* or *false* ( $U_T = U_F = 10$  MXN). Conversely, misclassifications, whether mistakenly classifying an accurate headline as *false* or a fake headline as *true*, yielded a utility of 0 MXN ( $U_{FT} = U_{TF} = 0$  MXN). The classification and WTP for each headline decision had a 20-second time limit. This utility structure incentivized participants to classify accurately.

For each headline, participants also indicated their willingness to pay (WTP) to access a perfect signal ( $S$ ) that could reveal the headline’s actual status. The perfect signal, available for purchase, would reveal the actual state of each headline with certainty ( $P(s = t|T) = q_t = 1$  and  $P(s = f|F) = q_f = 1$ ).

### Headline Number 18

Time left to complete this page: 0:01

Please classify the following headline: (If your classification is correct, you could earn an extra 10 MXN.)

**Iran Censored the Olympics; All Women Appear with Rectangles or Asterisks Covering Them**

Your Classification:

☐ The information is accurate ☒ Contains false information

How much are you willing to pay to verify this news?

1.5

Next

FIGURE 1. SCREENSHOT OF THE TRANSLATED CLASSIFICATION-VERIFICATION GAME AS SEEN BY THE PARTICIPANTS.

After completing all classifications and WTP decisions within a block, participants reported their estimated probability of correctly classifying the headlines by themselves and others. There was no time limit when participants reported their confidence. This provided a self-assessed confidence measure for the block.

After estimating their probabilities at the end of each block, participants received feedback according to the treatment group to which they were randomly assigned. Feedback treatments were designed to inform participants about their classification performance, either individually or relative to others. The control group was used as a reference. The feedback types and descriptions are presented in 1. By block, all treatment groups were shown a summary of the times they classified a headline as *true* or *false*, and the feedback treatments were shown the accuracy rates conditional on the headlines classified as *true* or *false*.

This structure allowed researchers to observe how different feedback types influenced participants’ accuracy, confidence, and valuation of the verification signal throughout the experiment.

TABLE 1—FEEDBACK TREATMENTS

Treatment Group	Feedback at the End of the Block
<b>Control Group</b>	No feedback on accuracy was given.
<b>Individual Feedback</b>	Personal accuracy rate for the block conditional on the headlines participants classified as <i>accurate</i> or <i>fake</i> .
<b>Others Feedback</b>	Average accuracy rate of other participants conditional on the headlines others classified as <i>accurate</i> or <i>fake</i> .

### C. Experimental Procedures

The participants were 184 undergraduate students in Mexico. The average age was 20 years old, and 55% of them were women. They were recruited from UNAM (National Autonomous University of Mexico) and IPN (National Polytechnic Institute), the first and second most important public schools in Mexico<sup>7</sup>. The experiment occurred at the schools where participants were studying in September 2024.

Participants receive a utility of 10 Mexican Pesos (MXN) for each correctly classified headline, whether it is accurate or fake ( $U_T = U_F = 10$  MXN). Conversely, they receive a utility of 0 MXN for misclassifications, whether they mistakenly classify an accurate headline as fake or a fake headline as accurate ( $U_{TF} = U_{FT} = 0$  MXN). This setup incentivizes participants to classify accurately and to value the signal appropriately based on its accuracy-assurance potential.

After all blocks, participants complete a survey collecting demographic data and assessing their support for the current government. For payment, one block is randomly selected at the experiment’s end, and participants receive 10 MXN for each correctly classified headline in the chosen block, thus linking final earnings directly to classification accuracy.

### D. Methodological Considerations

This experiment controls for variables that are relevant in the field to focus purely on the effects of feedback. The believed probability of receiving fake news, the interest in classifying correctly, and the verification quality change the value of verifying. The key parameters and assumptions—such as the equal prior probabilities, equal utilities, and a perfect signal—simplify the problem and allow for a focus on the probabilistic aspects of classification and verification.

<sup>7</sup>In the national ranking, UNAM is the most important university, and IPN can be ranked third (<https://www.usnews.com/education/best-global-universities/mexico>) or fourth (<https://www.topuniversities.com/university-rankings-articles/world-university-rankings/best-universities-mexico>), depending of the ranking.

We used the BDM mechanism to measure participants' WTP for verification. This was presented as a second price auction against a computer randomly choosing numbers from 0 to 5. Participants also chose a number between 0 and 5 from a drop-down menu to participate in the auction. If they win the auction, the signal is verified, and they win 10 MXN independently of the classification they made of the headline. If they lose the auction by betting less than the computer, their payoff will depend only on their initial classification.

In the field, verifying is a discrete decision based on each headline's expected gains and costs. However, the willingness to pay to verify directly measures the expected gains hidden in a discrete decision. Two people verifying (or not) can have different values for the information. Eliciting the WTP allows a continuous measurement of verification's value and, therefore, demand and how this might change depending on treatment.

We used the method in (Wilson and Vespa, 2018) to measure the confidence levels to present the binarized scoring rule in plain text. To increase this measure's reliability, we implemented this simple description of the problem (Charness et al., 2021), and assuring the participants that it is in their best interest to report their true beliefs (Danz et al., 2022) was implemented. The exact wording of this elicitation can be found in figure B1 of the appendix.

#### HEADLINES SELECTION

The online publication AnimalPolitico<sup>8</sup> and VerificadoMX<sup>9</sup> were used as the sources to find relevant fake news circulating in Mexico. These are the most relevant fact-checking efforts recognized in Mexico. To find headlines that were real but difficult to classify, the authors used NewsGPT<sup>10</sup>. The authors verified these headlines independently. All the headlines generated were verified independently by the authors. The authors selected 60 headlines from these sources: 30 political and 30 non-political, half of which were true and the other half false. Also, from the political headlines, 15 were classified as information that favored the government, and 15 opposed the government.

To select the 50 headlines for the final experiment and the order in which the blocks were placed, we ran a study in Prolific among Mexicans whose first language was Spanish. We asked for the classification of the headlines with the same incentives as in the final experiment and measured the probability of each headline being classified correctly. The headline composition of the blocks was made so that they have similar difficulty levels. The list of headlines used in the experiment can be found in the appendix table B.B2.

Using neutral headlines minimizes the potential influence of motivated reason-

<sup>8</sup><https://animalpolitico.com/verificacion-de-hechos>

<sup>9</sup><https://verificado.com.mx/>

<sup>10</sup>The request was made in August, around three weeks before the start of the first session: <https://chatgpt.com/g/g-NnU2wmnZ5-news-gpt-chat-with-hundreds-of-news-sources/c/7e750031-b534-481c-83cf-2dc6917d98b4>

ing, allowing the experiment to focus on the effects of feedback and overconfidence. At the same time, political news can show the impact of motivated reasoning on the demand for information.

#### POLITICAL POSITION

In the exit survey, participants had to answer some questions about politics and their answers to the questions: 1) who did you vote for? and 2) if the approved work made by AMLO<sup>11</sup> as president. If a participant answered "Morena"<sup>12</sup> to the first question and "Agree" to the second, that person was classified as a supporter of the government. If a participant voted for any other party and disagreed with the statement, they were classified as opposing the government. The polarization in Mexico has increased just as in the US. However, the main division is in options about AMLO and the political party he founded to run for president (Morena)<sup>13</sup>. This is an opinion shared in traditional media<sup>14</sup>, academics and journalists<sup>15</sup>. The experiment happened in September 2024, more than three months after the Presidential election in Mexico.

### V. Results

Table 2 shows the average values of the essential variables grouped by treatment. The first four variables are participants' characteristics; ages didn't change much between treatments, and the proportion of female participants was slightly larger in the Individual group, but in general, the groups were balanced. The proportion of participants supporting the government ranged from 16.4% to 24%, while the proportion of participants opposing the government ranged from 14.5%. Also, the proportion of missing headlines due to not answering within the 20-second limit was always less than 5%. Regarding confidence, their average probability of making the right decision is always higher than that of others. Also, their estimated accuracy rate is below the accuracy rate observed ("Correct"). The results in the accuracy and correctness suggest underestimation and overplacement (following Moore and Healy (2008)'s classification).

The effects on confidence, willingness to pay, and accuracy, can be observed in table 3 and 4. The regressions in table 3 were run at the level of blocks, while the regressions in table 4 at the level of round. This difference is because the confidence was measured at the end of a block of 10 rounds. Therefore, matching each round with the confidence level could be done by extrapolating

<sup>11</sup>This is a common way that people and media use to refer to Andres Manuel Lopez Obrador. The president of Mexico from 2018 to 2024.

<sup>12</sup>This is the name of the incumbent party during the federal elections in 2024.

<sup>13</sup><https://apnews.com/article/mexico-election-polarized-divided-heat-violence-4d5f620f0f8f9b7ef6efa8b3083561a8>

<sup>14</sup><https://apnews.com/article/mexico-election-polarized-divided-heat-violence-4d5f620f0f8f9b7ef6efa8b3083561a8>

<sup>15</sup><https://www.eluniversal.com.mx/tendencias/la-reflexion-de-denise-maerker-sobre-las-elecciones-2024-que-se-volvio-viral/>

TABLE 2—SUMMARY BY TREATMENT OF THE MAIN VARIABLES. THE AVERAGE (PROPORTION) OF EACH VARIABLE IS PRESENTED FOR EACH TREATMENT.

Variable	Control	Individual	Others
Age	19.9	20	20.4
Female	0.452	0.672	0.509
Support Gov	0.21	0.239	0.164
Oppose Gov	0.145	0.209	0.218
Missing Headlines	0.029	0.023	0.048
Accuracy Estimate	0.57	0.522	0.548
Accuracy Estimate Others	0.522	0.505	0.494
Accuracy	0.62	0.606	0.596
Classification ( $c = a$ )	0.497	0.507	0.511
WTP	2.82	2.64	2.42
N Participants	62	67	55

the confidence measure to all the headlines or taking the average of variables per round. There is incomplete information on the relationship between confidence, WTP, and accuracy. Here, we take the average as it is more conservative than extrapolating values. However, all the results here are robust enough to be present with the extrapolation approach. These results can be seen in the appendix. Also, the first block was dropped from this analysis because there was no feedback in that block, and the effects should have been observed after this block. The number of observations is also not multiples of 10 because some rounds were dropped out when the participants didn't provide an answer in the 20 seconds they had to.

#### A. Confidence

This section summarizes the analysis of the confidence participants had in their ability to classify headlines with false information. In regression 1 of table 3, the levels of confidence (measured by the reported  $\hat{P}(T|t)$  and  $\hat{P}(F|f)$ ) are explained by the treatments, participants characteristics, and headline properties. In this regression, the average WTP and Accuracy were taken so they are measured at the same level as the confidence elicited after each block of ten headlines.

After making their classifications, participants are asked to report the probability that they believe their classifications are correct when they classified a headline as true  $\hat{P}(T|t)$  and false  $\hat{P}(F|f)$ . This self-reported probability allows for the calculation of overconfidence metrics. Overconfidence is assessed by comparing the participants' reported probabilities of correct classification with the actual probabilities derived from the task. Specifically, overconfidence for accurate classifications is calculated as  $\hat{P}(T|t) - P(T|t)$ , and for fake classifications as  $\hat{P}(F|f) - P(F|f)$ . An overall measure of overconfidence is also calculated as

TABLE 3—REGRESSION ON THE CONFIDENCE AND WILLINGNESS TO PAY. THE REGRESSION IS ON THE DATA AT THE LEVEL OF BLOCKS TAKING THE AVERAGE OF THE ROUND VARIABLES. THE SE WERE CLUSTERED AT THE INDIVIDUAL LEVEL TO ACCOUNT FOR WITHIN-PARTICIPANT CORRELATION WHEN THERE ARE MULTIPLE OBSERVATIONS PER PARTICIPANT.

	<i>Dependent variable:</i>		
	Confidence	WTP	Accuracy
	(1)	(2)	(3)
Individual Feedback	−4.963 (3.029)	−0.190 (0.221)	−0.014 (0.016)
Others Feedback	−3.064 (3.198)	−0.374* (0.216)	−0.020 (0.016)
Block	−0.358 (0.749)	0.022 (0.030)	0.029*** (0.011)
’True’ (c = t)	1.063 (1.006)	0.205*** (0.042)	−0.007 (0.006)
Accuracy	−0.797 (2.517)	−0.399** (0.188)	
Age	0.932 (0.604)	−0.035 (0.054)	0.007** (0.003)
Male	3.156 (2.591)	0.007 (0.184)	0.010 (0.013)
Confidence		0.001 (0.003)	−0.0001 (0.0002)
Political	3.454** (1.521)	0.158** (0.065)	−0.010 (0.026)
Support Gov	7.974** (3.380)	0.494** (0.216)	0.041** (0.016)
Against Gov	5.047 (3.199)	0.275 (0.238)	0.0003 (0.016)
Constant	33.552*** (12.623)	3.326*** (1.091)	0.391*** (0.073)
Observations	1,464	1,464	1,464
R <sup>2</sup>	0.053	0.048	0.032
Adjusted R <sup>2</sup>	0.046	0.041	0.025

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01



$\hat{P}(c = \omega) - P(c = \omega)$ , where  $c$  is the classification and  $\omega$  is the true state of the world.

Table 2 shows that the average confidence was in line with observed accuracy rates, indicating a realistic self-assessment. This aligns with prior findings in similar tasks where participants are directly incentivized for accuracy and can adjust beliefs based on task difficulty. Considering the average overconfidence per participant, the hypothesis of no difference from zero is rejected ( $p$ -value < 0.001) with a difference of  $-7.15$ . The same result is found when analyzing the differences by initial classification:  $\hat{P}(F|f) - P(F|f) = -4.51$  and  $\hat{P}(T|t) - P(T|t) = -8.57$  (both with  $p$ -value < 0.001).

**RESULT 1:** *Contrary to hypothesis 1 overconfidence in headline classification, participants were not overconfident.*

Our results indicate **no statistically significant impact of feedback**—whether individual or group—on participants’ reported confidence in their classifications. This outcome suggests that feedback in this context did not affect participants’ self-evaluated accuracy, even though feedback types varied across conditions. This finding aligns with other studies showing limited feedback effects on confidence when feedback does not address specific decision outcomes.

Participants demonstrated higher confidence levels when classifying **political headlines**, with a notable increase among **government supporters**. This trend indicates that confidence is context-sensitive, with political relevance amplifying participants’ certainty in their classifications. Supporters of the government, in particular, tended to exhibit a stronger conviction that their classifications were accurate, especially for politically aligned headlines. However, having a political position against the government did not affect the level of confidence. We will analyze the political headlines in more detail in the second part of the results. These results are robust when a regression is considered at the headline level. These regressions can be seen in the appendix.

**RESULT 2:** *Confidence in headline classification increased when the headline was political and when they support the government.*

#### B. Willingness to Pay (WTP) for Verification

This section shows the results on the willingness to pay (WTP). In table 3, it can be observed that receiving feedback on **others’ accuracy feedback** lowered participants’ WTP to verify classifications in a statistically significant way. This reduction in demand for verification suggests that collective feedback may diminish the perceived need for individual verification, possibly due to an implicit assumption of greater task simplicity or shared accuracy. This outcome underscores the role of social feedback in shaping verification behavior in information classification tasks. This result can also be observed in figure 2 where the distribution of WTP for the control group second order stochastically dominates the distribution of WTP in the “Others” treatment.

RESULT 3: *Participants who received feedback on the performance of others were **less** willing to pay to verify headlines.*

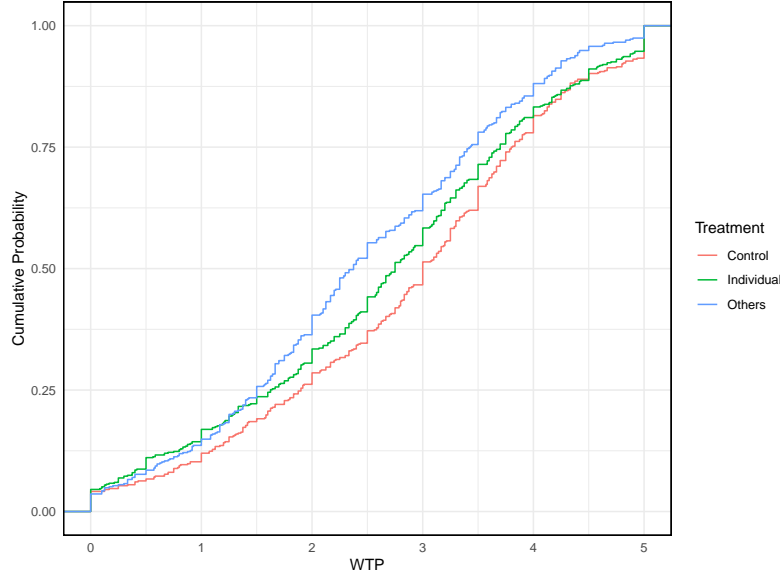


FIGURE 2. EMPIRICAL CDF OF THE WILLINGNESS TO PAY BY TREATMENT. TO CREATE THIS GRAPH, THE AVERAGE WTP PER BLOCK WAS CALCULATED.

Regression 2 of table 3 shows the significant negative effect of classifying a headline as *true*, and when their classification was larger. This indicates a verification bias, whereby participants seek to confirm rather than challenge initial beliefs. This finding is consistent with behavioral patterns in decision-making under uncertainty, where individuals are more likely to invest in information that reinforces their prior beliefs. This result can also be observed in figure 3 where the distribution of WTP for the *false* classification second order stochastically dominates the distribution of WTP in *true* classification. This is true for every treatment.

RESULT 4: *Participants exhibited a **greater** willingness to pay (WTP) to verify the information they classified as true.*

Participants demonstrated an **increased WTP for verifying political news**, emphasizing the perceived importance of accuracy for politically sensitive content. This effect was amplified among **government supporters**, who displayed a higher demand for verification, potentially reflecting a higher stake in the perceived accuracy of politically favorable information. This finding suggests that

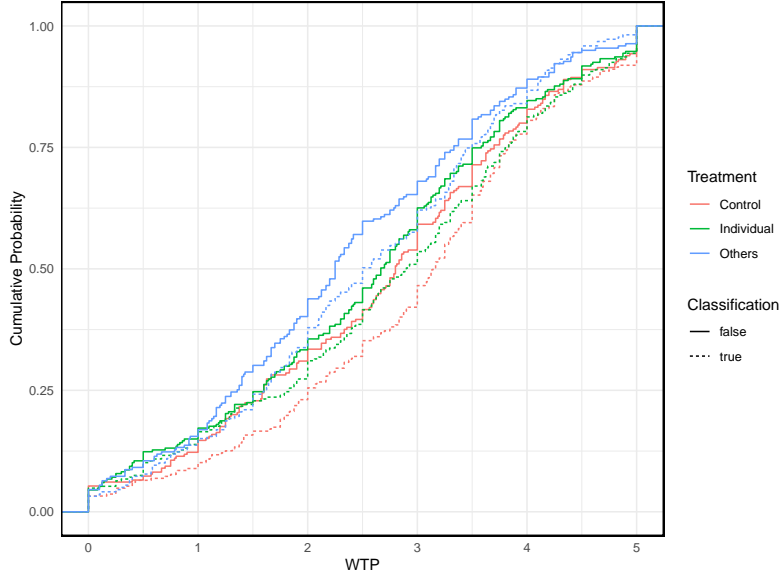


FIGURE 3. EMPIRICAL CDF OF THE WILLINGNESS TO PAY BY TREATMENT AND CLASSIFICATION OF THE HEADLINE. TO CREATE THIS GRAPH, THE AVERAGE WTP PER BLOCK WAS CALCULATED.

political alignment influences classification behavior and the demand for verification.

### C. Effects on Political Headlines

In this section, we analyze only the headlines containing political content that can be classified as "Favorable" or "Unfavorable" to the Government. The results can be observed in the regressions in table 4.

The WTP results observed across all headlines were preserved when analyzing political news exclusively. However, in contrast to general findings, **political alignment did not significantly affect WTP** within this subset, indicating that demand for verification in politically charged content may be driven by factors other than simple partisan alignment. Also, consistent with results from all the headlines, participants were **more likely** to classify a headline as accurate if it favored the government, a trend especially pronounced among those who opposed the government. This suggests that partisan perception plays a role in classification tendencies, with individuals more likely to label favorable headlines as accurate even when they hold opposing political views. Additionally, for headlines correctly classified initially, participants were more inclined to classify them as accurate, reinforcing the tendency to affirm initial judgments. Also, the effects of feedback on others' accuracy were significant.

TABLE 4—REGRESSION ON THE WILLINGNESS TO PAY AND CLASSIFICATION. THE REGRESSION IS ON THE DATA AT THE LEVEL OF ROUNDS EXTRAPOLATING THE BLOCK CONFIDENCE TO THE CONFIDENCE IN EACH HEADLINE. THE SE WERE CLUSTERED AT THE INDIVIDUAL LEVEL TO ACCOUNT FOR WITHIN-PARTICIPANT CORRELATION WHEN THERE ARE MULTIPLE OBSERVATIONS PER PARTICIPANT.

	<i>Dependent variable:</i>	
	WTP	Accuracy
	(1)	(2)
Individual Feedback	−0.214 (0.240)	−0.019 (0.017)
Others Feedback	−0.450* (0.241)	−0.040** (0.019)
Block	0.042 (0.050)	−0.137*** (0.016)
Confidence	0.001 (0.004)	0.00003 (0.0003)
'True' (c = t)	0.242*** (0.072)	0.207*** (0.006)
Accuracy	−0.040 (0.054)	
Age	−0.048 (0.056)	0.005 (0.004)
Male	0.057 (0.202)	0.011 (0.015)
Support Gov	0.392 (0.255)	0.016 (0.023)
News Favor Gov	−0.033 (0.045)	−0.216*** (0.016)
Against Gov	0.254 (0.261)	−0.063*** (0.024)
Support Gov X Favor Gov	0.092 (0.082)	−0.002 (0.033)
Favor Gov X Against Gov	0.083 (0.079)	0.091*** (0.031)
Constant	3.420*** (1.160)	1.181*** (0.104)
Observations	3,603	3,603
R <sup>2</sup>	0.033	0.088
Adjusted R <sup>2</sup>	0.029	0.085

*Note:*

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

Interestingly, there was a negative effect of being in the Others treatment on the task accuracy, only when looking at the political headlines. This means that people who knew the performance of others did performed worse when classifying political news. This effect was absent when considering all the headlines. Also, there was a robust negative effect of the block. It seems unlikely that this is an effect of learning as this is captured by the block, and it is more likely that knowing the others' accuracy has a larger impact on the attention set on the task when the headlines are political. Due to the nature of politics, it is possible that people pay more attention to information that they perceive as more socially relevant.

In regression 2 of table 4, it can be seen that the accuracy was larger when the classification was *true*. This result was also found when analyzing all the blocks in table 3. There were no effects of the interactions between the political position of the participants and the political leaning of the headline on the WTP. This refutes hypothesis 3.

RESULT 5: *Accuracy is **larger** when participants classify a headline as true.*

In terms of the effects of political variables, if the headlines favor the government and participants are against it, the accuracy of the headline classification is reduced. There is also a positive interaction between the headline favoring the government and the participant being against it. If people who criticize the government are better informed of what is false when the information favors the the government.

## VI. Discussion

We found a negative relationship between the feedback on other's performance and the willingness to pay for information. Considering that underconfidence was a more prevalent trait of the participants, this result follows the theory; they could be learning that the task is easier than expected, and they will consider that the probability of correctly classifying is higher, which reduces the value of new verifying the headline. However, we didn't find a significant effect of feedback on the probability they thought they got a correct classification (level of confidence). The willingness to pay was asked by each headline (50 observations), while the confidence was asked at the end of the block. This mismatch creates a noisy relation between these variables since the average confidence is not linked to individual headlines, where participants can be very confident or very unsure about the veracity of the headline. We believe that this noise in the measure of both variables hides the relationship between the expected probability of correct classification and willingness to pay to verify.

### A. Perfect Verification

The methodology presented here simplifies the world, allowing us to explore the causal effects of feedback on the value of verification. In the field, we have

verification practices that could be imperfect and different utilities for believing fake news and rejecting accurate information. Providing feedback in the real world could change the importance of each kind of mistake. This could be behind the effects of highlighting the importance of accuracy observed in [Pennycook et al. \(2021\)](#) and [Pennycook and Rand \(2022\)](#).

The best way to verify is an open question not addressed in the current research; the question is about overconfidence as a mechanism behind low levels of verification. We present a perfect signal to avoid motivated misinterpretation of the signal’s likelihoods ([Thaler \(2024\)](#)) and to avoid the discussion of the real information value of a specific verification practice. The results from [Thaler \(2024\)](#) are significant since a no-completely-informative signal will allow a biased updating belief that also decreases the WTP for the signal. Even online searching of news, one of the most common practices promoted in digital literacy programs, could backfire ([Aslett et al. \(2024\)](#); [Hoes et al. \(2023\)](#)). Also, having imperfect signals might make purchasing any information suboptimal if participants expect their original classification to remain the same even after the

### B. Classification and Verification Behavior

The results indicate that participants classified headlines as *true* ( $c = t$ ) approximately 50% of the time ( $P(c = t) = 0.5$ ), aligning with the known prior probabilities of each state being equal ( $P(T) = P(F) = 0.5$ ), which they were informed about. In terms of confidence, there was a slight asymmetry:  $P(A|c = t) = 54.7\%$  and  $P(F|c = f) = 53.6\%$ . This difference was minimal, suggesting that participants may adhere to the martingale property, as their classification accuracy remained consistent with the prior.

Interestingly, participants showed a greater willingness to pay (WTP) for verification when they classified a headline as  $(c = t)$ , suggesting that they were more motivated to confirm their accurate classifications. This finding implies an inequality that can be expressed through the decision-making analysis:

$$\begin{aligned}
 WTP(c = t) &> WTP(c = f) \\
 &\Leftrightarrow \\
 1 - P(T|c = t) &> 1 - P(F|c = f) \\
 &\Leftrightarrow \\
 P(A|c = t) &< P(F|c = f)
 \end{aligned}$$

Additionally, participants were more accurate when classifying headlines as accurate:  $P(c = T|c = t) = 67.3\%$  compared to  $P(c = F|c = f) = 60.2\%$ . This contradicts the general expectation that people would exhibit more caution in marking information as accurate. One possible explanation is that participants

may have an internalized utility structure that penalizes certain types of misclassifications, such as inaccurately labeling false information as accurate (e.g.,  $U_{TF} < 0$ ). However, if such a utility structure were in effect, we might expect a lower overall proportion of headlines classified as accurate due to increased caution.

These findings highlight subtle asymmetries in confidence and verification demand that could inform future research on verification behavior, especially regarding how individuals internalize error costs in classification tasks.

### C. Incentives

The incentives were relevant for the sample. Participants were paid 10 Mexican Pesos (MXN) per correctly classified headline from a block selected at random. This is approximately 0.5 USD per headline. This means that they could make between 0 and 100 MXN from the classification-verification game. In Mexico the average salary of a profesionist is 53.5 MXN.<sup>16</sup> Then, considering that this task required less than an hour (including instructions and a survey) and that the participants were undergrad students, the money they received was enough. If they decide not to pay attention to the task, they could have maximized their expected payoffs by selecting a willingness to pay of 5 MXN. Since this was a second price auction against a random process, the expected value would be 75 MXN. However, not many rounds were selected with a WTP of 5 MXN. For the Control treatment, it was 18.02 %; for the Individual treatment, it was 13.23 %; and for the other treatment, it was 10.2 %. These results coincide with the general finding of a reduced demand for verification in the Others treatment. This might signal that participants infer that the task was easy after the feedback. Therefore, they prefer to verify the headlines themselves and reduce the verification amount.

### D. Time

In this section we present the most important results in the time people took in each round.<sup>17</sup> Participants had a limit of 20 seconds to classify the headlines and indicate their willingness to pay. As noted in table 2, the proportion of missing headlines was always below 5 %. And this is the reason the number of observations in the regression is not a multiple of 5.

In figure 4, we can see that there was no difference in the distribution of time spend on headlines that where classified as *true* and *false*. Also the mode is at 10 seconds. It is clear that 20 seconds were enough for most participants in most of the headlines and that if there were participants waiting until last second and potentially making wrong decisions, they represent a small proportion of the rounds.

<sup>16</sup>Data from the Mexican government.

<sup>17</sup>The experiment was deployed in Heroku where the time spend in each page has registered.

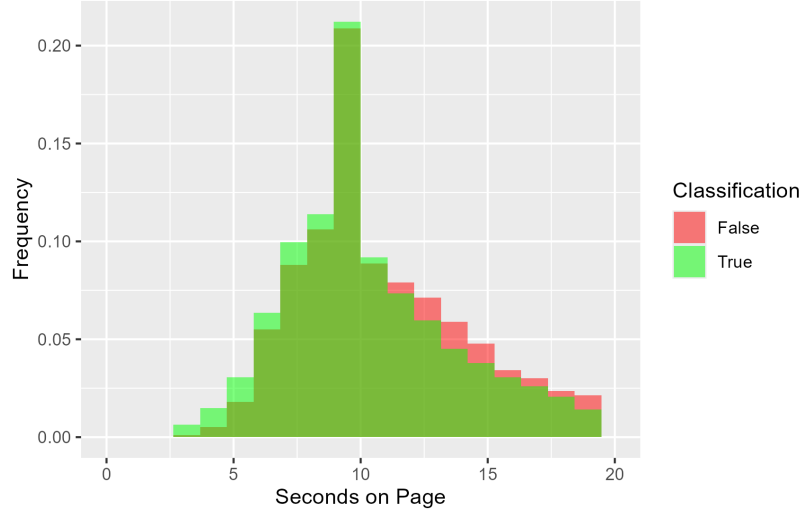


FIGURE 4. FREQUENCY OF TIME SPEND ON HEADLINES THAT WHERE CLASSIFIED AS TRUE AND FALSE.

The headlines were also difficult enough so that the proportion of wrong classifications is not zero, but too difficult such that the participants choose randomly. In figure 2, we can see that 'real' headlines (that were actually True) have larger proportion of classifications as *true*. As a reference, the red line indicates an even classification among participants, which means classifying the headline was difficult. Although some headlines are close to this line (between the blue lines in 0.4 and 0.6) most headlines were well classified by the participants. We expected that difficult headline took more time, and there is a group of headlines with this behavior, however most headlines were classified in around 10 seconds as also noted in figure 4. Other analysis of the time can be found in the appendix.

## VII. Conclusion

This study provides novel insights into how feedback influences individuals' willingness to pay (WTP) for verifying news headlines, their confidence in classification, and their accuracy in identifying misinformation. By incorporating feedback mechanisms in a controlled experimental setting, we examined how different feedback types shaped decision-making processes and their implications for combating misinformation.

One of the most significant findings is that feedback on the accuracy of others significantly reduced WTP for verification, supporting Hypothesis 2. Participants in the "Others" feedback condition demonstrated a lower demand for verification compared to those receiving individual or no feedback. This reduction suggests that social feedback leads participants to infer that the classification task is easier than initially perceived, diminishing the perceived need for verification. However,



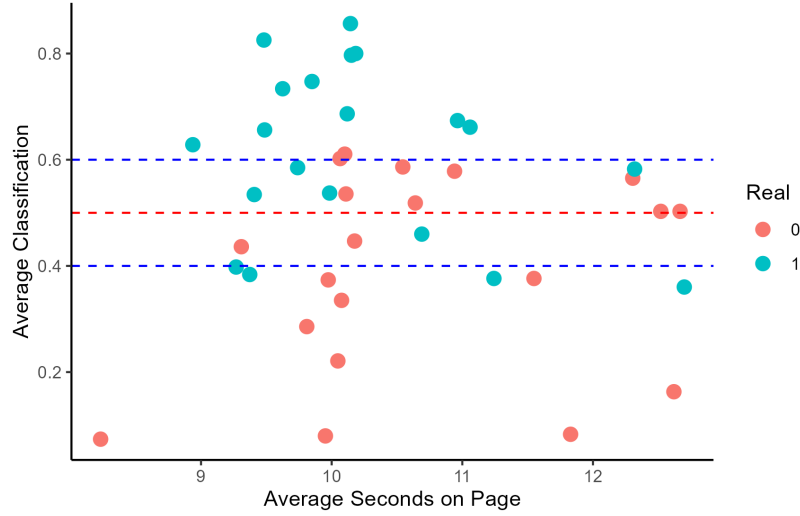


FIGURE 5. AVERAGE TIME SPENT ON EACH HEADLINE AND THE PROPORTION OF CORRECT CLASSIFICATIONS AGAINST THE AVERAGE CLASSIFICATION AS "TRUE".

contrary to expectations, feedback did not significantly affect self-reported confidence levels, suggesting that while participants adjusted their verification behaviors, their confidence remained stable. This decoupling between confidence and verification highlights the nuanced relationship between perceived task difficulty and the value of additional information.

When examining political content, participants exhibited a higher WTP for verifying politically charged headlines, particularly those aligned with their political views. However, we did not find significant interactions between participants' political alignment and the leaning of the headlines on WTP, refuting Hypothesis 3. This suggests that while political content affects verification demand, it is not directly moderated by the alignment between participants' political positions and the headlines' content. Nonetheless, participants who supported the government were more confident when classifying headlines favorable to their political stance, demonstrating how political alignment amplifies perceived accuracy. These results underscore the role of political bias in shaping both confidence and verification behaviors.

Interestingly, participants exhibited a verification bias, being more willing to pay to verify headlines they initially classified as *true*, compared to those classified as *false*. This result aligns with Hypothesis 1, which predicted that confidence would influence WTP for verification. However, instead of overconfidence leading to lower verification demand, as hypothesized, participants appeared to seek confirmation for their classifications, revealing an asymmetry in how they valued information. This suggests that interventions targeting verification behaviors

should account for individuals’ tendencies to confirm, rather than challenge, their initial judgments.

Lastly, while personal feedback improved classification accuracy compared to no feedback, as proposed in Hypothesis 4, the effect was modest and primarily observed for politically neutral content. Furthermore, participants were generally more accurate when classifying headlines as *true*, reinforcing the notion that individuals are more cautious in affirming rather than rejecting information. Interestingly, participants who received “Others” feedback performed worse when classifying political headlines, possibly due to divided attention between task difficulty and the social relevance of political content. These findings highlight the importance of understanding the contexts in which feedback can enhance or hinder classification accuracy.

#### A. *Implications and Future Research*

Our findings suggest that feedback interventions, particularly those involving peer performance metrics, need careful calibration to avoid unintended decreases in verification demand. The lack of overconfidence observed in our experiment implies that individuals may already possess a reasonably accurate self-assessment of their information classification skills in structured settings. However, the observed confirmation bias in verification demand calls for further exploration into how individuals’ motivations for verification might be nudged to encourage unbiased fact-checking.

These results reveal the complex interplay between feedback, confidence, verification demand, and political motivations. While feedback on others’ accuracy reduced WTP for verification, it did not consistently improve confidence or accuracy. This finding has critical implications for designing interventions to combat misinformation: feedback mechanisms must consider the cognitive processes underlying verification demand and confidence. Additionally, the influence of political content on verification behavior underscores the importance of addressing motivated reasoning and political bias in combating misinformation.

Future research should explore the dynamics of imperfect verification signals and the role of feedback in shaping long-term verification behaviors. By expanding the experimental framework to include real-world verification mechanisms and varying levels of signal accuracy, we can better understand how to design interventions that effectively reduce misinformation in diverse informational environments. This study provides a foundation for further exploration of how feedback and political biases shape decision-making in the context of misinformation, with practical implications for policy, education, and platform design.

## REFERENCES

- Arin, K Peren, Deni Mazrekaj, and Marcel Thum**, “Ability of detecting and willingness to share fake news,” *Scientific Reports*, 2023, 7298.
- Arrow, K. J., D. Blackwell, and M. A. Girshick**, “Bayes and Minimax Solutions of Sequential Decision Problems,” *Econometrica*, 7 1949, 17, 213.
- Aslett, Kevin, Zeve Sanderson, William Godel, Nathaniel Persily, Jonathan Nagler, and Joshua A. Tucker**, “Online searches to evaluate misinformation can increase its perceived veracity,” *Nature*, 1 2024, 625, 548–556.
- Bateman, Jon and Dean Jackson**, *Countering Disinformation Effectively: An Evidence-Based Policy Guide*, Carnegie Endowment for International Peace, 2024.
- Charness, Gary, Uri Gneezy, and Vlastimil Rasocha**, “Experimental methods: Eliciting beliefs,” *Journal of Economic Behavior and Organization*, 2021, 189, 234–256.
- Danz, David, Lise Vesterlund, and Alistair J Wilson**, “Belief Elicitation and Behavioral Incentive Compatibility,” *American Economic Review*, 2022, 112, 2851–2883.
- Eberlein, Marion, Sandra Ludwig, and Julia Nafziger**, “The effects of feedback on self-assessment,” *Bulletin of Economic Research*, 2011, 63, 307–3378.
- Erat, Serhat, Kurtulus Demirkol, and M Eyyüp Sallabas**, “Overconfidence and its link with feedback,” *Active Learning in Higher Education*, 2022, 3, 173–187.
- Ferraro, Paul J**, “Know thyself: incompetence and overconfidence,” *Experimental Laboratory Working Paper Series No. 2003-001. Department of Economics, Andrew Young School of Policy Studies, Georgia State University*, 2005.
- Hoes, Emma, Brian Aitken, Jingwen Zhang, Tomasz Gackowski, and Magdalena Wojcieszak**, “Prominent misinformation interventions reduce misperceptions but increase scepticism,” *Nature Human Behavior*, 2023.
- Kartal, Melis and Jean Robert Tyran**, “Fake News, Voter Overconfidence, and the Quality of Democratic Choice,” *American Economic Review*, 10 2022, 112, 3367–97.
- Kogelnik, Maria**, “Performance Feedback and Gender Differences in Persistence,” *SSRN Electronic Journal*, 12 2022.

- Kozyreva, Anastasia, Philipp Lorenz-Spreen, Stefan M. Herzog, Ullrich K.H. Ecker, Stephan Lewandowsky, Ralph Hertwig, Ayesha Ali, Joe Bak-Coleman, Sarit Barzilai, Melisa Basol, Adam J. Berinsky, Cornelia Betsch, John Cook, Lisa K. Fazio, Michael Geers, Andrew M. Guess, Haifeng Huang, Horacio Larreguy, Rakoén Maertens, Folco Panizza, Gordon Pennycook, David G. Rand, Steve Rathje, Jason Reifler, Philipp Schmid, Mark Smith, Briony Swire-Thompson, Paula Szewach, Sander van der Linden, and Sam Wineburg, “Toolbox of individual-level interventions against online misinformation,” *Nature Human Behaviour* 2024 8:6, 5 2024, 8, 1044–1052.
- Lyons, Benjamin A., Jacob M. Montgomery, Andrew M. Guess, Brendan Nyhan, and Jason Reifler, “Overconfidence in news judgments is associated with false news susceptibility,” *Proceedings of the National Academy of Sciences of the United States of America*, 6 2021, 118.
- Moore, Don A. and Paul J. Healy, “The Trouble With Overconfidence,” *Psychological Review*, 4 2008, 115, 502–517.
- Oprea, Ryan and Sevgi Yuksel, “Social Exchange of Motivated Beliefs,” *Journal of the European Economic Association*, 2022, 20, 667–699.
- Ortoleva, Pietro and Erik Snowberg, “Overconfidence in Political Behavior,” *American Economic Review*, 2015, 105, 504–535.
- Pennycook, Gordon and David G. Rand, “Who falls for fake news? The roles of bullshit receptivity, overclaiming, familiarity, and analytic thinking,” *Journal of Personality*, 4 2020, 88, 185–200.
- and —, “Nudging Social Media toward Accuracy,” *Annals of the American Academy of Political and Social Science*, 3 2022, 700, 152–164.
- , Ziv Epstein, Mohsen Mosleh, Antonio A Arechar, Dean Eckles, and David G Rand, “Shifting attention to accuracy can reduce misinformation online,” *Nature* 592, 2021, 592, 590–595.
- Pulford, Briony D. and Andrew M. Colman, “Overconfidence: Feedback and item difficulty effects,” *Personality and Individual Differences*, 7 1997, 23, 125–133.
- Thaler, Michael, “The Fake News Effect: Experimentally Identifying Motivated Reasoning Using Trust in News,” *American Economic Journal: Microeconomics*, 2024, 16, 1–38.
- Vosoughi, Soroush, Deb Roy, and Sinan Aral, “The spread of true and false news online,” *Science*, 2018, 359, 1146–1151.
- Wald, Abraham, “Foundations of a General Theory of Sequential Decision Functions,” *Econometrica*, 10 1947, 15, 279.

**Wilson, Alistair J and Emanuel Vespa**, “Paired-uniform scoring: Implementing a binarized scoring rule with non-mathematical language,” 2018.

## MATHEMATICAL APPENDIX

## A1. Expected Value of the Signal Considering Reclassification

## GENERAL SETUP

- 1) **Initial Classification:** The agent initially classifies the headline as  $c$  (either accurate ( $t$ ) or fake ( $f$ )).
- 2) **Receive Signal:** The agent receives a signal  $s$  which can either confirm or contradict their initial classification.
- 3) **Reclassification:** Based on the signal, the agent makes a new classification  $c'$ .

EXPECTED UTILITY WITH SIGNAL AND RECLASSIFICATION ( $EU_{\text{SIGNAL}}^{\text{UPDATE}}(c)$ )

The expected utility with the signal and reclassification is calculated by considering the updated posterior probabilities and the new classification based on the signal.

## STEP 1: DEFINE PROBABILITIES AND UTILITIES

• **Initial Posterior Probabilities:**

$$P(T|c) = \frac{P(c|T) \cdot (1 - p_f)}{P(c)}, \quad P(F|c) = \frac{P(c|F) \cdot p_f}{P(c)}$$

where:

$$P(c) = (1 - p_f) \cdot P(c|T) + p_f \cdot P(c|F)$$

• **Signal Probabilities:**

$$q_t = P(s = t|T), \quad q_f = P(s = f|F)$$

• **Utilities:**

$$U_T, U_F, U_{TF}, U_{FT}$$

## STEP 2: DEFINE UPDATED POSTERIOR PROBABILITIES GIVEN SIGNAL

After observing the signal, the agent updates their beliefs:

• **Posterior Probabilities Given Signal  $s = t$ :**

$$P(T|s = t, c) = \frac{q_t \cdot P(T|c)}{q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)}$$

$$P(F|s = t, c) = \frac{(1 - q_f) \cdot P(F|c)}{q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)}$$

• **Posterior Probabilities Given Signal  $s = f$ :**

$$P(T|s = f, c) = \frac{(1 - q_t) \cdot P(T|c)}{(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)}$$

$$P(F|s = f, c) = \frac{q_f \cdot P(F|c)}{(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)}$$

STEP 3: CALCULATE EXPECTED UTILITY AFTER SIGNAL

We assume that the signal is strong enough (proposition 1) to assure that the best reclassification is to follow what the signal indicates is the state of the world. Otherwise, the signal would have no instrumental value, and therefore,  $EVSI = 0$ .

The expected utility with the signal, considering the possibility of reclassification, is:

$$EU_{\text{signal}}^{\text{update}}(c) = P(s = t|c) \cdot EU_{\text{new classification}}(s = t, c) + P(s = f|c) \cdot EU_{\text{new classification}}(s = f, c)$$

Here,  $P(s = t|c)$  and  $P(s = f|c)$  are the probabilities of receiving the signals  $s = t$  and  $s = f$  given the initial classification  $c$ . These probabilities are determined by Bayes' rule, considering the agent's initial classification and the properties of the signal.

Given the initial classification  $c$ :

$$P(s = t|c) = q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)$$

$$P(s = f|c) = (1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)$$

EXPECTED UTILITY AFTER SIGNAL  $s = t$

$$EU_{\text{new classification}}(s = t, c) = P(T|s = t, c) \cdot U_T + P(F|s = t, c) \cdot U_{TF}$$

EXPECTED UTILITY AFTER SIGNAL  $s = f$

$$EU_{\text{new classification}}(s = f, c) = P(F|s = f, c) \cdot U_F + P(T|s = f, c) \cdot U_{FT}$$

## STEP 4: COMBINE EXPECTED UTILITIES

$$EU_{\text{signal}}^{\text{update}}(c) = [q_t \cdot P(T|c) + (1 - q_f) \cdot P(F|c)] \cdot [P(T|s = t, c) \cdot U_T + P(F|s = t, c) \cdot U_{TF}] \\ + [(1 - q_t) \cdot P(T|c) + q_f \cdot P(F|c)] \cdot [P(F|s = f, c) \cdot U_F + P(T|s = f, c) \cdot U_{FT}]$$

## STEP 5: EXPECTED VALUE OF THE SIGNAL (EVSI)

Finally, the EVSI is the difference between the expected utility with the signal and the expected utility without the signal.

$$EVSI = EU_{\text{signal}}^{\text{update}}(c) - EU_{\text{no signal}}(c)$$

## CONCLUSION

By allowing the agent to update their classification based on the signal, we account for the dynamic decision-making process. The expected value of the signal (EVSI) is derived by comparing the expected utility with the signal (considering reclassification) to the expected utility without the signal. This approach shows the impact of additional information on improving decision-making accuracy.

*A2. Proof of Proposition 2: Need for a Strong Enough Signal*

To ensure that people follow the signal  $S$  for reclassification, we must prove that the expected utility of reclassifying based on the signal is higher than not reclassifying. Without loss of generality, we will first consider the case where the initial classification  $c = f$  and the signal  $s = t$ .

## INITIAL SETUP

- 1) **Initial Classification:**  $c = f$  (classified as fake)
- 2) **Signal Received:**  $s = t$  (signal indicates true)

We need to show that reclassifying the headline as true ( $c' = t$ ) based on the signal is optimal.

## EXPECTED UTILITY OF NOT RECLASSIFYING

If the agent does not reclassify and sticks with the initial classification  $c = f$ , but knows the signal  $s = t$ , the expected utility is:

$$EU_{\text{no reclassification}}(f, s = t) = P(T|s = t, f) \cdot U_{FT} + P(F|s = t, f) \cdot U_F$$



## EXPECTED UTILITY OF RECLASSIFYING

If the agent reclassifies the headline based on the signal  $s = t$ , the expected utility is:

$$EU_{\text{reclassification}}(f, s = t) = P(T|s = t, f) \cdot U_T + P(F|s = t, f) \cdot U_{TF}$$

## POSTERIOR PROBABILITIES

The posterior probabilities given the signal  $s = t$  and initial classification  $c = f$  are:

$$P(T|s = t, f) = \frac{q_t \cdot P(T|f)}{q_t \cdot P(T|f) + (1 - q_f) \cdot P(F|f)}$$

$$P(F|s = t, f) = \frac{(1 - q_f) \cdot P(F|f)}{q_t \cdot P(T|f) + (1 - q_f) \cdot P(F|f)}$$

## CONDITION FOR RECLASSIFYING

To prove that reclassifying based on the signal is optimal, we need:

$$EU_{\text{reclassification}}(f, s = t) > EU_{\text{no reclassification}}(f, s = t)$$

Substituting the utilities, we get:

$$P(T|s = t, f) \cdot U_T + P(F|s = t, f) \cdot U_{TF} > P(T|s = t, f) \cdot U_{FT} + P(F|s = t, f) \cdot U_F$$

Given the simplifying assumptions:

$$U_T = 1, \quad U_F = 1, \quad U_{TF} = 0, \quad U_{FT} = 0$$

The inequality simplifies to:

$$P(T|s = t, f) \cdot 1 + P(F|s = t, f) \cdot 0 > P(T|s = t, f) \cdot 0 + P(F|s = t, f) \cdot 1$$

This reduces to:

$$P(T|s = t, f) > P(F|s = t, f)$$

## VERIFYING THE POSTERIOR PROBABILITIES

Substitute the posterior probabilities:

$$\frac{q_t \cdot P(T|f)}{q_t \cdot P(T|f) + (1 - q_f) \cdot P(F|f)} > \frac{(1 - q_f) \cdot P(F|f)}{q_t \cdot P(T|f) + (1 - q_f) \cdot P(F|f)}$$

Since the denominators are the same, we can simplify this to:

$$q_t \cdot P(T|f) > (1 - q_f) \cdot P(F|f)$$

Since  $P(F|f) = 1 - P(T|f)$ , we have:

$$q_t \cdot P(T|f) > (1 - q_f) \cdot (1 - P(T|f))$$

Expanding and rearranging terms, we get:

$$q_t \cdot P(T|f) > (1 - q_f) - (1 - q_f) \cdot P(T|f)$$

$$q_t \cdot P(T|f) + (1 - q_f) \cdot P(T|f) > (1 - q_f)$$

$$P(T|f) \cdot (q_t + 1 - q_f) > (1 - q_f)$$

Dividing both sides by  $(q_t + 1 - q_f)$ :

$$P(T|f) > \frac{1 - q_f}{q_t + 1 - q_f}$$

This shows that the signal needs to be strong enough such that  $q_t$  is sufficiently large compared to  $1 - q_f$ , ensuring that the agent reclassifies the headline as true based on the signal. This proves that a strong signal is necessary to ensure that people follow the signal  $S$  for reclassification.

TRIVIAL CASE:  $c = s = t$

If the initial classification  $c = t$  and the signal  $s = t$ , then reclassification is not necessary because the initial classification is already true. The expected utility remains the same:

$$EU_{\text{reclassification}}(t, s = t) = P(T|s = t, t) \cdot U_T + P(F|s = t, t) \cdot U_{TF}$$

Given the simplifying assumptions, this reduces to:

$$EU_{\text{reclassification}}(t, s = t) = P(T|s = t, t) \cdot 1 + P(F|s = t, t) \cdot 0 = P(T|s = t, t)$$

The expected utility of not reclassifying is:

$$EU_{\text{no reclassification}}(t, s = t) = P(T|s = t, t) \cdot U_T + P(F|s = t, t) \cdot U_{TF}$$

Given the simplifying assumptions, this reduces to:

$$EU_{\text{no reclassification}}(t, s = t) = P(T|s = t, t) \cdot 1 + P(F|s = t, t) \cdot 0 = P(T|s = t, t)$$

Since both expected utilities are equal, reclassification is trivial in this case.

## OTHER CASES

The same process follows for the cases  $c = t, s = f$  and  $c = s = f$ . For these cases, the conditions are as follows:

1) **Case**  $c = t, s = f$ :

$$EU_{\text{reclassification}}(t, s = f) > EU_{\text{no reclassification}}(t, s = f)$$

Substituting the utilities, we get:

$$P(F|s = f, t) \cdot U_F + P(T|s = f, t) \cdot U_{FT} > P(F|s = f, t) \cdot U_{TF} + P(T|s = f, t) \cdot U_T$$

Given the simplifying assumptions:

$$P(F|s = f, t) \cdot 1 + P(T|s = f, t) \cdot 0 > P(F|s = f, t) \cdot 0 + P(T|s = f, t) \cdot 1$$

This reduces to:

$$P(F|s = f, t) > P(T|s = f, t)$$

Verifying the posterior probabilities:

$$\frac{q_f \cdot P(F|t)}{q_f \cdot P(F|t) + (1 - q_t) \cdot P(T|t)} > \frac{(1 - q_t) \cdot P(T|t)}{q_f \cdot P(F|t) + (1 - q_t) \cdot P(T|t)}$$

Since the denominators are the same, we can simplify this to:

$$q_f \cdot P(F|t) > (1 - q_t) \cdot P(T|t)$$

Since  $P(T|t) = 1 - P(F|t)$ , we have:

$$q_f \cdot P(F|t) > (1 - q_t) \cdot (1 - P(F|t))$$

Expanding and rearranging terms, we get:

$$q_f \cdot P(F|t) > (1 - q_t) - (1 - q_t) \cdot P(F|t)$$

$$q_f \cdot P(F|t) + (1 - q_t) \cdot P(F|t) > (1 - q_t)$$

$$P(F|t) \cdot (q_f + 1 - q_t) > (1 - q_t)$$

Dividing both sides by  $(q_f + 1 - q_t)$ :

$$P(F|t) > \frac{1 - q_t}{q_f + 1 - q_t}$$

2) **Case**  $c = s = f$ : If the initial classification  $c = f$  and the signal  $s = f$ , then reclassification is unnecessary because the initial classification is already

correct. The expected utility remains the same:

$$EU_{\text{reclassification}}(f, s = f) = P(F|s = f, f) \cdot U_F + P(T|s = f, f) \cdot U_{TF}$$

Given the simplifying assumptions, this reduces to:

$$EU_{\text{reclassification}}(f, s = f) = P(F|s = f, f) \cdot 1 + P(T|s = f, f) \cdot 0 = P(F|s = f, f)$$

The expected utility of not reclassifying is:

$$EU_{\text{no reclassification}}(f, s = f) = P(F|s = f, f) \cdot U_F + P(T|s = f, f) \cdot U_{TF}$$

Given the simplifying assumptions, this reduces to:

$$EU_{\text{no reclassification}}(f, s = f) = P(F|s = f, f) \cdot 1 + P(T|s = f, f) \cdot 0 = P(F|s = f, f)$$

Since both expected utilities are equal, reclassification is trivial in this case.

#### CONDITIONS

Therefore, to ensure that the signal is strong enough to prompt optimal reclassification in both cases, we need to satisfy two key conditions:

$$P(T|f) > \frac{1 - q_f}{q_t + 1 - q_f}$$

$$P(T|t) < \frac{q_f}{q_f + 1 - q_t}$$

Considering the conditions for proposition 1 we have that,

$$\frac{1 - q_f}{q_t + 1 - q_f} < P(T|f) < P(T|t) < \frac{q_f}{q_f + 1 - q_t}$$

*A3. Proof of Proposition 1: Sufficient and Necessary Condition for*

$$P(T|c = f) < P(T) < P(T|c = t) \text{ and } P(F|c = t) < P(F) < P(F|c = f)$$

*A4. Proof of Necessity and Sufficiency*

We will prove that  $1 < \frac{P(c=t|T)}{P(c=t|F)}$  and  $1 < \frac{P(c=f|F)}{P(c=f|T)}$  if and only if  $P(T|c = f) < P(T) < P(T|c = t)$  and  $P(F|c = t) < P(F) < P(F|c = f)$ .

#### DEFINITIONS AND SETUP

Let:

- $P(T)$  be the prior probability that the state is true.

- $P(F)$  be the prior probability that the state is fake.
- $P(c = t|T)$  be the probability of classifying a headline as true given it is true.
- $P(c = t|F)$  be the probability of classifying a headline as true given it is fake.
- $P(c = f|F)$  be the probability of classifying a headline as fake given it is fake.
- $P(c = f|T)$  be the probability of classifying a headline as fake, given it is true.

#### POSTERIOR PROBABILITIES

The posterior probabilities after observing the classification  $c$  are given by:

- Posterior probability of  $T$  given  $c = f$ :

$$P(T|c = f) = \frac{P(c = f|T) \cdot P(T)}{P(c = f|T) \cdot P(T) + P(c = f|F) \cdot P(F)}$$

- Posterior probability of  $T$  given  $c = t$ :

$$P(T|c = t) = \frac{P(c = t|T) \cdot P(T)}{P(c = t|T) \cdot P(T) + P(c = t|F) \cdot P(F)}$$

- Posterior probability of  $F$  given  $c = f$ :

$$P(F|c = f) = \frac{P(c = f|F) \cdot P(F)}{P(c = f|T) \cdot P(T) + P(c = f|F) \cdot P(F)}$$

- Posterior probability of  $F$  given  $c = t$ :

$$P(F|c = t) = \frac{P(c = t|F) \cdot P(F)}{P(c = t|T) \cdot P(T) + P(c = t|F) \cdot P(F)}$$

#### PART 1: SUFFICIENCY ( $\Rightarrow$ )

Assume that  $1 < \frac{P(c=t|T)}{P(c=t|F)}$  and  $1 < \frac{P(c=f|F)}{P(c=f|T)}$ . We want to show that this implies  $P(T|c = f) < P(T) < P(T|c = t)$  and  $P(F|c = t) < P(F) < P(F|c = f)$ .

## ANALYZE THE POSTERIOR PROBABILITIES

- 1) **For**  $P(T|c = f)$ : Given the condition  $1 < \frac{P(c=f|F)}{P(c=f|T)}$ , we know that:

$$\frac{P(c = f|F)}{P(c = f|T)} > 1$$

This implies  $P(c = f|F) > P(c = f|T)$ . As a result, in the posterior probability expression:

$$P(T|c = f) = \frac{P(c = f|T) \cdot P(T)}{P(c = f|T) \cdot P(T) + P(c = f|F) \cdot P(F)}$$

The denominator  $P(c = f|T) \cdot P(T) + P(c = f|F) \cdot P(F)$  will be larger than the numerator  $P(c = f|T) \cdot P(T)$ , causing  $P(T|c = f)$  to be smaller than the prior  $P(T)$ . Therefore:

$$P(T|c = f) < P(T)$$

- 2) **For**  $P(T|c = t)$ : Given the condition  $1 < \frac{P(c=t|T)}{P(c=t|F)}$ , we know that:

$$\frac{P(c = t|T)}{P(c = t|F)} > 1$$

This implies  $P(c = t|T) > P(c = t|F)$ . As a result, in the posterior probability expression:

$$P(T|c = t) = \frac{P(c = t|T) \cdot P(T)}{P(c = t|T) \cdot P(T) + P(c = t|F) \cdot P(F)}$$

The numerator  $P(c = t|T) \cdot P(T)$  will dominate the denominator  $P(c = t|T) \cdot P(T) + P(c = t|F) \cdot P(F)$ , causing  $P(T|c = t)$  to be larger than the prior  $P(T)$ . Therefore:

$$P(T|c = t) > P(T)$$

- 3) **For**  $P(F|c = f)$  **and**  $P(F|c = t)$ : Similarly, the same reasoning applies to  $P(F|c = f)$  and  $P(F|c = t)$ , given that:

$$\frac{P(c = f|F)}{P(c = f|T)} > 1 \quad \text{and} \quad \frac{P(c = t|T)}{P(c = t|F)} > 1$$

This implies that:

$$P(F|c = t) < P(F) < P(F|c = f)$$

PART 2: NECESSITY ( $\Leftarrow$ )

Assume that  $P(T|c = f) < P(T) < P(T|c = t)$  and  $P(F|c = t) < P(F) < P(F|c = f)$ . We need to show that this implies  $1 < \frac{P(c=t|T)}{P(c=t|F)}$  and  $1 < \frac{P(c=f|F)}{P(c=f|T)}$ .

## ANALYZING THE POSTERIOR PROBABILITIES

- \*\*For  $P(T|c = f) < P(T)$ :\*\*

Given the posterior probability expression:

$$P(T|c = f) = \frac{P(c = f|T) \cdot P(T)}{P(c = f|T) \cdot P(T) + P(c = f|F) \cdot P(F)}$$

If  $P(T|c = f) < P(T)$ , then the likelihood ratio  $\frac{P(c=f|F)}{P(c=f|T)}$  must be greater than 1. This is because the posterior  $P(T|c = f)$  being less than  $P(T)$  implies that the signal  $c = f$  is more likely to come from the fake state  $F$ , meaning:

$$\frac{P(c = f|F)}{P(c = f|T)} > 1$$

- \*\*For  $P(T|c = t) > P(T)$ :\*\*

Given the posterior probability expression:

$$P(T|c = t) = \frac{P(c = t|T) \cdot P(T)}{P(c = t|T) \cdot P(T) + P(c = t|F) \cdot P(F)}$$

If  $P(T|c = t) > P(T)$ , then the likelihood ratio  $\frac{P(c=t|T)}{P(c=t|F)}$  must be greater than 1. This is because the posterior  $P(T|c = t)$  being greater than  $P(T)$  implies that the signal  $c = t$  is more likely to come from the true state  $T$ , meaning:

$$\frac{P(c = t|T)}{P(c = t|F)} > 1$$

- \*\*For  $P(F|c = f) > P(F)$  and  $P(F|c = t) < P(F)$ :\*\*

By symmetry, the same reasoning applies for  $P(F|c = f) > P(F)$  and  $P(F|c = t) < P(F)$ . The likelihood ratios  $\frac{P(c=f|F)}{P(c=f|T)} > 1$  and  $\frac{P(c=t|T)}{P(c=t|F)} > 1$  are necessary conditions to satisfy these posterior inequalities.

## CONCLUSION

Thus, we have shown that:  $1 < \frac{P(c=t|T)}{P(c=t|F)}$  and  $1 < \frac{P(c=f|F)}{P(c=f|T)}$  are necessary and sufficient conditions for:

$$P(T|c=f) < P(T) < P(T|c=t)$$

$$P(F|c=t) < P(F) < P(F|c=f)$$

## COROLLARY: INFORMATIVENESS OF THE SIGNAL

The same analysis can be applied to the signal. Therefore  $1 < \frac{P(s=t|T)}{P(s=t|F)}$  and  $1 < \frac{P(s=f|F)}{P(s=f|T)} \iff$

$$P(T|s=f) < P(T) < P(T|s=t)$$

$$P(F|s=t) < P(F) < P(F|s=f)$$



## MATERIALS

*B1. Confidence Elicitation***Confidence in Block Classification 3**

Answer the following questions with the probability in percentage terms.  
Where 100 means the event always occurs, 0 means it never occurs, and 50 means it occurs half of the time.

**Please consider the block of 10 news headlines that you just classified:**

You classified 5 headlines as "The information is accurate" and 5 as "Contains false information".

One of the 5 headlines you classified as accurate will be selected at random.  
What is the probability that the headline is actually accurate?

One of the 5 headlines you classified as false will be selected at random.  
What is the probability that the headline is actually false?

**Now, consider the classification that other participants made in this block of 10 news headlines:**

A headline classified as accurate by another participant will be selected at random.  
What is the probability that the headline is actually accurate?

A headline classified as false by another participant will be selected at random.  
What is the probability that the headline is actually false?



FIGURE B1. SCREENSHOT OF THE TRANSLATED CONFIDENCE ELICITATION AS SEEN BY THE PARTICIPANTS.

*B2. Headlines Used in the Experiment*

Block	Real	Headline	Translated Headline
1	1	Se inaugura un nuevo museo en honor a Cantinflas en la Ciudad de México	A New Museum in Honor of Cantinflas is Inaugurated in Mexico City
1	1	El salario mínimo en México se incrementa 20% en 2024	Minimum Wage in Mexico Increases by 20% in 2024
1	1	La variante Ómicron es la única de preocupación que circula a nivel mundial; es más transmisible, aunque menos peligrosa que la variante Delta	The Omicron Variant is the Only Variant of Concern Circulating Worldwide; It Is More Transmissible, Though Less Dangerous than the Delta Variant

1	1	Se suspende programa humanitario para trabajar o solicitar asilo en Estados Unidos para Haití, Venezuela, Nicaragua y Cuba	Humanitarian program to work or apply for asylum in the United States for Haiti, Venezuela, Nicaragua, and Cuba is suspended
1	1	Incrementó en el uso de energías renovables en México	Increase in the Use of Renewable Energy in Mexico
1	0	Turismo internacional se desploma en 2024, México ya no es un destino atractivo	International Tourism Collapses, Mexico is No Longer an Attractive Destination
1	0	Luto en México por accidente aéreo de un avión de pasajeros. No hubo sobrevivientes	National Mourning in Mexico. Terrible Passenger Plane Crash in 2024. No Survivors
1	1	En 2024, se intensificaron los incendios forestales en México	In 2024, Wildfires Intensified in Various Regions of Mexico
1	1	Existen programas de apoyo a pequeñas empresas lanzados por el gobierno mexicano	There Are Support Programs for Small Businesses Launched by the Mexican Government
1	0	Se firma un tratado del Foro Económico Mundial que busca reconocer la pedofilia como orientación sexual	A World Economic Forum Treaty Is Signed to Recognize Pedophilia as a Sexual Orientation
2	1	El INAH cobrará \$60 por tomar fotografías para uso comercial en museos y sitios arqueológicos	INAH Will Charge \$60 for Taking Photos in Museums and Archaeological Sites for Commercial Use
2	0	Se publica la lista de apellidos que pueden solicitar la ciudadanía española	Spain Publishes a List of Surnames That Allow One to Apply for Spanish Citizenship
2	0	En Irán censuraron los Juegos Olímpicos; todas las mujeres aparecen con rectángulos o asteriscos cubriéndolas	Iran Censored the Olympics; All Women Appear with Rectangles or Asterisks Covering Them
2	0	Iniciará juicio en contra de la ministra presidenta de la SCJN por participar en el paro de trabajadores del Poder Judicial.	Trial Against the Chief Justice of the Supreme Court for Participating in the Judicial Workers' Strike to Begin
2	0	La Organización de Estados Americanos (OEA) sanciona a México por dar asilo a Jorge Glass en la embajada mexicana en Ecuador	The Organization of American States (OAS) Managed to Sanction Mexico for Granting Asylum to Jorge Glass in the Mexican Embassy
2	0	El hijo de Nicolás Maduro es capturado en video manejando un Ferrari dorado	Nicolás Maduro's son is seen driving a golden Ferrari

2	1	El Ozempic, promovido en redes para bajar de peso, es un tratamiento controlado para la diabetes tipo 2	Ozempic Is Actually a Controlled Treatment for Type 2 Diabetes
2	1	México alcanza cifra récord en exportaciones agrícolas	Mexico Reaches Record High in Agricultural Exports
2	0	Atletas ucranianos portaron pulseras de tobillo con GPS para evitar su huida después de los Juegos Olímpicos de París 2024	Ukrainian Athletes Wore GPS Ankle Bracelets to Prevent Them from Fleeing After the Olympic Games
2	1	Ningún país ha declarado confinamiento por mpox tras la nueva emergencia sanitaria anunciada por la OMS	No country has declared a lockdown due to mpox following the new health emergency announced by the WHO
3	0	La cama "antisexo" siguió siendo utilizada en los Juegos Olímpicos de París 2024	The "Anti-Sex" Bed Will Continue to Be Used at the Paris 2024 Olympics
3	0	El aeropuerto de Suecia fue descontaminado debido a contagios de Mpox	Sweden's Airport Was Decontaminated Due to Mpox Infections
3	0	Miss Venezuela protestó ante las cámaras contra su gobierno en una alfombra roja	Miss Venezuela Protested Against Her Government on a Red Carpet
3	1	Avances en la investigación de nuevas vacunas desarrolladas en México	Advances in the Research of New Vaccines Developed in Mexico
3	0	Un estadounidense se suicidó saltando desde su habitación durante el Baja Beach Fest 2024 en México	An American Committed Suicide During the Baja Beach Fest 2024 in Mexico
3	1	Descubrimiento de nuevas ruinas mayas en la península de Yucatán en 2024	Discovery of New Mayan Ruins in the Yucatán Peninsula
3	1	México termina en el puesto 65 del medallero en los Juegos Olímpicos de París 2024	Mexico Finishes 65th in the Medal Table at the Paris 2024 Olympic Games
3	1	Mexico envió dos aviones en 2023 para rescatar connacionales varados en Israel por el conflicto en Gaza	Sedena and SRE Sent Two Planes to Rescue Mexicans Stranded in Israel
3	0	Consumir alimentos alcalinos ayuda a contrarrestar la variante Omicron del coronavirus	Maintaining a pH (Acidity Level) Above 5.5 Can Prevent Covid-19 Infection

3	0	El Consejo para Prevenir y Eliminar la Discriminación (COPRED) busca suspender la celebración del Día del Padre en los centros educativos	COPRED Urges Elementary Schools Not to Exclude Children from Non-Normative Families on Father's Day
4	1	México tiene la tasa más baja de desempleo de la OCDE	Mexico Has the Lowest Unemployment Rate in the OECD
4	0	Gobierno de México entrega el nuevo "Bono Mujeres" por 2 mil 700 pesos	Mexican Government Issues the New "Women's Bonus" for 2,700 Pesos
4	1	Se registró una disminución de 5.1 millones personas en pobreza en el actual gobierno	A Decrease of 5.1 Million People in Poverty Was Recorded During the Current Government
4	0	Tribunal Electoral encuentra irregularidades graves en el triunfo de Claudia Sheinbaum	Electoral Tribunal Finds Serious Irregularities in Claudia Sheinbaum's Victory
4	1	La presidenta electa Claudia Sheinbaum anuncia beca universal para estudiantes de nivel básico	President-Elect Claudia Sheinbaum Announces Universal Scholarship for Elementary School Students
4	0	El Tren Maya se completará sin impacto ambiental, según estudios científicos independientes	The Maya Train Will Be Completed Without Environmental Impact, According to Independent Scientific Studies
4	1	11 Ministros de la Suprema Corte de Justicia de la Nación ganan \$206,246 pesos mensuales netos	11 Supreme Court Justices Earn \$206,246 Pesos Monthly Net
4	1	México envió dos aviones a Israel para rescatar a las y los mexicanos varados por el conflicto con Palestina.	Mexico Sent Two Planes to Israel to Rescue Mexicans Stranded Due to the Conflict with Palestine
4	1	En solo ocho de cada 100 delitos en México se abre una carpeta de investigación	Only Eight Out of Every 100 Crimes in Mexico Lead to an Investigation Being Opened
4	0	México prepara una reunión con los presidentes de Rusia y Corea del Norte para comprar armas	Mexico Prepares a Meeting with the Presidents of Russia and North Korea to Buy Weapons
5	1	En el gobierno de AMLO se logró una reducción en la tasa de homicidios.	AMLO's Government Achieved a Reduction in the Homicide Rate
5	1	Primer director femenino de la CFE es nombrado en México	First Female Director of the CFE Is Appointed in Mexico
5	0	Durante el gobierno de Andrés Manuel López Obrador, la deuda pública subió 64%	During Andrés Manuel López Obrador's Government, Public Debt Increased by 64%

5	0	Poder Judicial de la Federación hay 53,737 personas que ganan más que el Presidente	In the Federal Judiciary, 53,737 People Earn More Than the President
5	1	La pobreza extrema en México incrementó de 2018 a 2022	Extreme Poverty in Mexico Increased from 2018 to 2022
5	1	Hay déficit presupuestario en el 2024 por parte del gobierno de México	Mexico's Government Faces a Budget Deficit in 2024
5	1	EU critica al Gobierno de AMLO por 'desacreditar a periodistas'	U.S. Criticizes AMLO's Government for 'Discrediting Journalists'
5	1	Sedena gastó más que el presupuesto autorizado por el Congreso	Sedena Spent More Than the Budget Authorized by Congress
5	0	México ya produce el 90% de la gasolina que consume, como afirma AMLO	Mexico Now Produces 90% of the Gasoline It Consumes, As Claimed by AMLO
5	0	Metro de CDMX dejará de ser gratis para adultos mayores	CDMX Metro Will No Longer Be Free for Senior Citizens

## REGRESSIONS APPENDIX

TABLE C1—REGRESSION ON THE CONFIDENCE AND WILLINGNESS TO PAY. THE REGRESSION IS ON THE DATA AT THE LEVEL OF ROUNDS EXTRAPOLATING THE BLOCK CONFIDENCE TO THE CONFIDENCE IN EACH HEADLINE. THE SE WERE CLUSTERED AT THE INDIVIDUAL LEVEL TO ACCOUNT FOR WITHIN-PARTICIPANT CORRELATION WHEN THERE ARE MULTIPLE OBSERVATIONS PER PARTICIPANT.

	<i>Dependent variable:</i>	
	Confidence	WTP
	(1)	(2)
Individual Feedback	−5.626* (3.045)	−0.196 (0.222)
Others Feedback	−3.119 (3.232)	−0.393* (0.218)
Block	−0.893 (0.817)	0.027 (0.029)
'Accurate' (c = a)	1.286 (1.101)	0.249*** (0.061)
Correct	−0.235 (0.406)	−0.059 (0.042)
Age	0.924 (0.641)	−0.037 (0.054)
Male	3.060 (2.649)	0.027 (0.184)
Confidence		0.001 (0.003)
Political	5.215*** (1.641)	0.130** (0.063)
Support Gov	7.763** (3.392)	0.471** (0.218)
Against Gov	4.793 (3.201)	0.261 (0.241)
Constant	35.120*** (12.778)	3.128*** (1.091)
Observations	7,173	7,173
R <sup>2</sup>	0.055	0.034
Adjusted R <sup>2</sup>	0.054	0.032

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

TABLE C2—REGRESSION ON THE WILLINGNESS TO PAY AND ACCURACY. THE REGRESSION IS ON THE DATA AT THE BLOCK LEVEL TAKING THE AVERAGE OF THE VARIABLES IN EACH BLOCK. THE SE WERE CLUSTERED AT THE INDIVIDUAL LEVEL TO ACCOUNT FOR WITHIN-PARTICIPANT CORRELATION WHEN THERE ARE MULTIPLE OBSERVATIONS PER PARTICIPANT.

	<i>Dependent variable:</i>	
	WTP	Accuracy
	(1)	(2)
Individual Feedback	−0.227*** (0.017)	−0.017 (0.017)
Others Feedback	−0.487*** (0.019)	−0.033* (0.019)
Block	0.043*** (0.015)	−0.140*** (0.015)
Confidence	0.001*** (0.0003)	0.00003 (0.0003)
'True' (c = t)	0.231*** (0.006)	0.210*** (0.006)
Accuracy	−0.039***	
WTP		−0.003 (0.004)
Age	−0.026*** (0.001)	0.0004 (0.001)
Male	0.045*** (0.015)	0.013 (0.015)
Support Gov	0.887*** (0.104)	−0.041 (0.104)
News Favor Gov	0.394*** (0.047)	−0.016 (0.047)
Against Gov	0.397*** (0.023)	0.017 (0.023)
Support Gov X Favor Gov	−0.029* (0.016)	−0.223*** (0.016)
Favor Gov X Against Gov	0.252*** (0.023)	−0.065*** (0.023)
support_govTRUE:Favor_gov	0.060* (0.032)	0.006 (0.032)
Favor_gov:against_gov	0.078** (0.031)	0.097*** (0.031)
Constant	2.998*** (0.082)	1.285*** (0.082)
Observations	3,762	3,762
R <sup>2</sup>	0.037	0.090
Adjusted R <sup>2</sup>	0.033	0.087

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

## C1. Time on Pages Figures

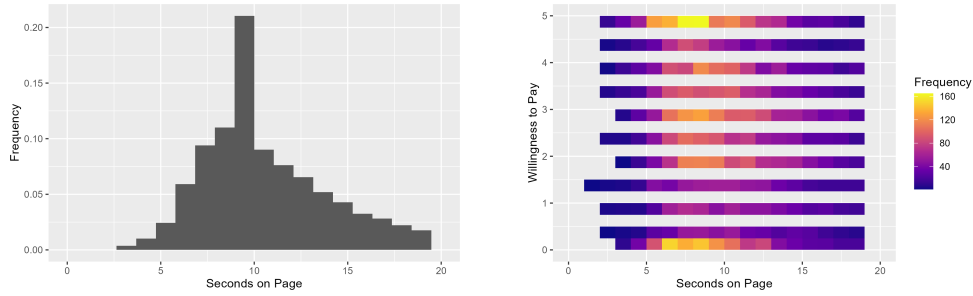


FIGURE C1. TIME SPEND ON EACH HEADLINE. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.

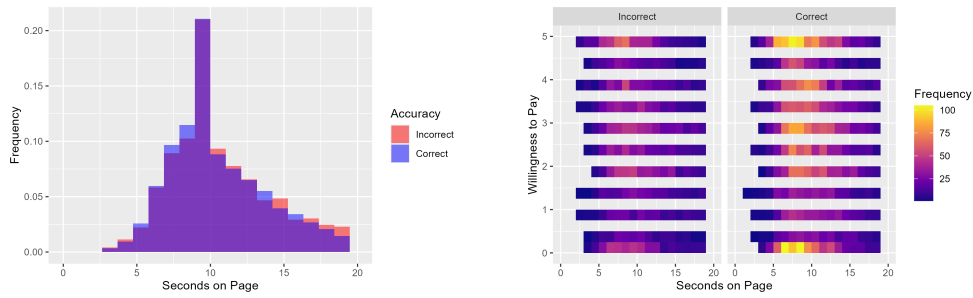


FIGURE C2. TIME SPEND ON EACH HEADLINE BY THE ACCURACY OF THE CLASSIFICATION. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.



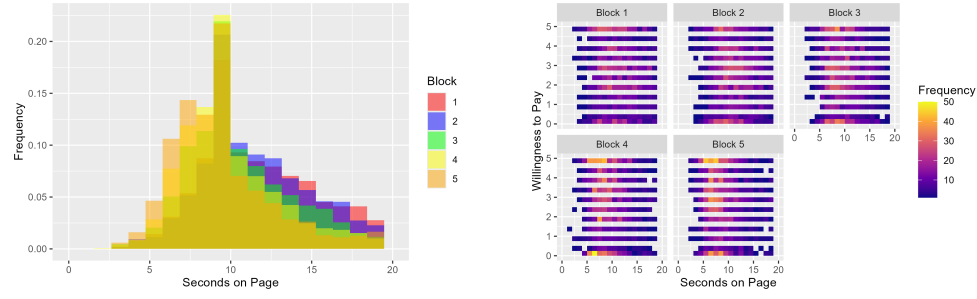


FIGURE C3. TIME SPEND ON EACH HEADLINE BY BLOCK. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.

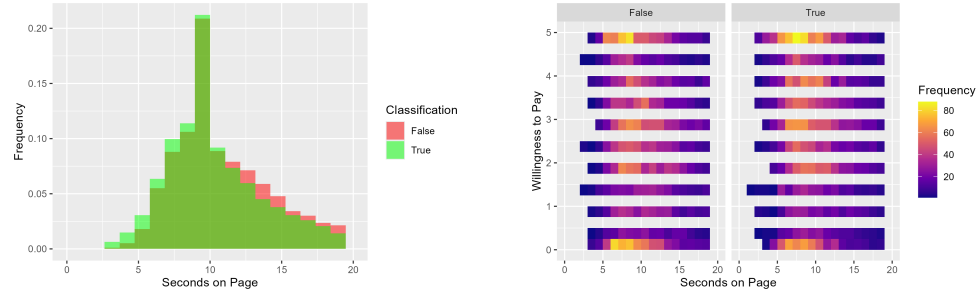


FIGURE C4. TIME SPEND ON EACH HEADLINE BY CLASSIFICATION. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.

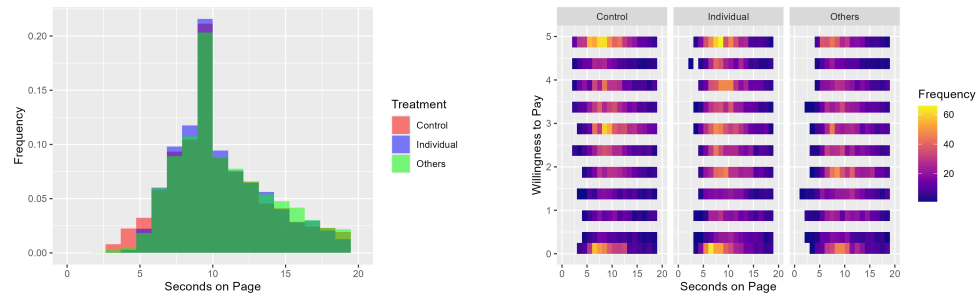


FIGURE C5. TIME SPEND ON EACH HEADLINE BY TREATMENT. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.

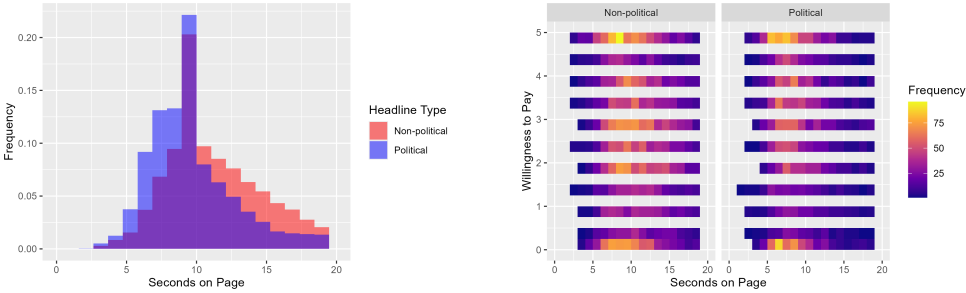


FIGURE C6. TIME SPEND ON EACH HEADLINE BY POLITICAL CONTENT. IN THE LEFT SIDE THERE IS THE WHOLE DISTRIBUTION, AND IN THE RIGHT SIDE THE DISTRIBUTIONS BY WILLINGNESS TO PAY LEVEL.