

Algorithmic Transparency and Portfolio Choices: Field Evidence

Béatrice Boulu-Reshef^a

Alexis Direr^b

Mehdi Louafi^c

January 8, 2026

Abstract

This paper studies whether profile-based explanations influence investors' acceptance of algorithmic risk recommendations in a randomized controlled trial embedded directly in the platform's interface of a leading French robo-advisor. Users were assigned either to see graphical explanations of the drivers underlying their recommended risk score and associated portfolio or to receive the standard interface with no explanation. Our results, obtained in a real-world setting with actual clients of a FinTech, do not support the adherence gains from increased transparency that are widely anticipated in the literature. Overall, providing profile-based explanations is not found to increase acceptance of the recommended profile nor raise users' engagement with the platform. However, we find a heterogeneous treatment effect as profile-based explanations lead to a greater downward deviation among desktop users who have already deviated to safer-than-recommended portfolios, but this pattern disappears once users' experience of the platform is taken into account. We observe non-causal evidence in both conditions that behavior is shaped primarily by the digital context and experience: phone and first-time users are more likely to accept the portfolio recommendation than desktop and returning users. While such transparency-enhancing profile-based explanations are informative, they are not a universal lever for adherence, suggesting that explanation design should be tested and tailored across device types and users' experience.

Keywords: Robo-Advisor, Financial advice, Portfolio Choices, Household Finance, Algorithmic Transparency

JEL classification: G50; G11; G23; C93

^aCY Cergy Paris Université, THEMA, CNRS, beatrice.boulu-reshef@cyu.fr

^bUniversité d'Orléans, LÉO, alexis.direr@univ-orleans.fr

^cUniversité d'Orléans, LÉO, mehdi.louafi@univ-orleans.fr

1 Introduction

Robo-advisors —algorithmic, largely automated portfolio advice delivered via digital platforms— promise scalable, low-cost guidance (Abraham, Schmukler, & Tessada, 2019; D’Acunto & Rossi, 2023), consistent with broader FinTech-driven declines in intermediation costs (Philippon, 2019), yet adoption remains modest (Cardillo & Chiappini, 2024). A frequently cited friction is algorithmic transparency: recommendations are perceived as “black-boxes,” users lack visibility into why a given risk profile fits their circumstances, and the loss of human interaction can erode trust (Morana, Gnewuch, Jung, & Granig, 2020; Jung, Dorner, Weinhardt, & Puzmaz, 2018). Explainable AI (XAI) proposes user-facing explanations to address this gap; yet, there is limited field evidence on whether explanations actually change behavior in high-stakes financial decisions (Adadi & Berrada, 2018; Anjomshoae, Najjar, Calvaresi, & Främling, 2019; Arrieta et al., 2020; Riedl, 2019; Weitz, Schiller, Schlagowski, Huber, & André, 2021). We therefore ask whether profile-based explanations, concise disclosures of the drivers of a recommended risk score benchmarked to similar users, increase acceptance of algorithmic risk profile recommendations on a market-deployed robo-advisor.

We study this question with a randomized controlled trial embedded in the interface of a leading French robo-advisor. The experiment randomized 4,645 saving-plan contracts to one of two conditions: (i) a treatment in which users received graphical, profile-based explanations detailing the relative importance of variables (economic situation, liquidity needs, investment horizon, demographics, risk tolerance) and benchmarking them to a “typical similar user”; or (ii) a control without any additional information beyond the recommended risk profile. Our primary outcome is acceptance, defined as choosing the recommended risk profile at the decision point. We also track the exact deviation from the recommended profile and engagement metrics. Sessions are observed on desktop and mobile, allowing us to study how the usage context moderates responses.

Three results stand out. First, explanations do not raise acceptance: treatment and control accept at statistically indistinguishable rates. Second, device context matters: desktop sessions exhibit lower baseline adherence than mobile. Third, among desktop sessions that choose a safer-than-recommended portfolio, explanations are associated with larger downward moves; however, this pattern loses significance once platform experience is accounted for. Beyond this subgroup, we find little systematic heterogeneity and no improvement in engagement with the platform. Taken together, profile-based transparency is not a universal lever for adherence, and its effects are context-dependent.

Our paper makes three contributions. First, we provide in-market experimental evidence on a widely advocated transparency intervention, documenting a precise null on acceptance and on deviations. Second, we show that behavior is shaped by both interface and experience: device context (desktop vs. mobile) and platform tenure (first-time vs. returning users) are first-order correlates of adherence—mobile and first-time users accept more, while desktop and returning users explore and deviate more; moreover, the desktop-specific amplification of downward moves is sensitive to controlling for experience. Third, we offer managerial guidance: for deployments seeking higher acceptance, adding profile-based explanations does not guarantee improved adher-

ence; rather, firms should test context-contingent designs (e.g., lighter explanations on desktop at the acceptance step, progressive disclosure, alternative benchmarks) instead of expecting uniform gains from transparency.

We implement a single-feature intervention by mirroring how a human advisor would justify a recommendation, grouping drivers into meaningful categories, and situating the user in a peer reference class, while preserving the product’s decision flow. The randomization was implemented within the production system and analyzed independently by the authors.

Our findings point to a clear null: in this randomized, in-market setting, profile-based explanations do not measurably change acceptance, the direction or magnitude of deviations, or engagement. Apparent interactions in simpler models are not robust once we account for user experience. The systematic patterns we observe, higher adherence on mobile and among first-time users, are driven by device and experience with the platform, not by the explanation itself.

The remainder of the paper proceeds as follows. Section 2 situates our study within research on explainability, robo-advising, and digital finance. Section 3 details the experimental design and measurement. Section 4 reports the main effects and robustness. Section 5 concludes with implications for the design of explainable robo-advice.

2 Related work

2.1 Algorithmic transparency and explanations

Concerns about the opacity of modern machine learning (ML) and Artificial Intelligence (AI) systems have fueled calls for algorithmic transparency from stakeholders, regulators,¹ and users (Castelvecchi, 2016; Preece, Harborne, Braines, Tomsett, & Chakraborty, 2018). XAI attempts to make system outputs more comprehensible to humans (Adadi & Berrada, 2018), with two complementary aims: (i) clarifying the drivers of individual or aggregate predictions (Anjomshoae et al., 2019; Arrieta et al., 2020; Beaudouin et al., 2020) and (ii) supporting understanding, acceptance, and trust (Cheng et al., 2019; Cai, Jongejan, & Holbrook, 2019; Shin, 2021; Weitz et al., 2021).

On the methods side, XAI distinguishes inherently interpretable models from black boxes equipped with post-hoc explanations (Arrieta et al., 2020). Explanations can be local, targeting a single decision (e.g., LIME (Ribeiro, Singh, & Guestrin, 2016)), or global, describing overall feature influence (e.g., SHAP/SAGE (Lundberg & Lee, 2017; Cohen, Dror, & Ruppin, 2007; Covert, Lundberg, & Lee, 2021); see also Speith (2022); Krzyziński, Spytek, Baniecki, and Biecek (2023)). A common design is feature-importance attribution and short rationales (“recommended because of A and B; despite C”), widely used in recommender systems (Nunes & Jannach, 2017).

On the human side, experiments in Human-Computer Interaction and Computer Science show that explanations typically improve understanding, but their effects on confidence and acceptance vary with context and presentation: example-based or agent-mediated explanations can help; yet, gains in comprehension do not always translate into greater trust or adoption (Cai et al., 2019;

¹Regulation (EU) 2016/679 (GDPR).

Shin, 2021; Weitz et al., 2021; Rader, Cotter, & Cho, 2018; Eslami, Krishna Kumaran, Sandvig, & Karahalios, 2018; Herlocker, Konstan, & Riedl, 2000; Cramer et al., 2008). Human-centered work stresses that “good” explanations must anticipate edge cases and align with user goals (Riedl, 2019). In sum, the literature offers mature tools for transparency but mixed behavioral predictions, motivating field evidence in high-stakes settings.

2.2 XAI, explanations, and robo-advisors

In robo-advising, opacity has been found to impact trust and delegation (Patel & Lincoln, 2019; Bianchi & Brière, 2020). Users often display algorithm aversion, thus preferring human judgment even when algorithms perform better (Dietvorst, Simmons, & Massey, 2015, 2018). Lab and online studies suggest that transparent explanations can mitigate aversion and support adoption; however, effects depend on format and timing. Accuracy or feature-based disclosures can raise uptake and preserve trust after errors (Ben David, Resheff, & Tron, 2021). Similar insights were found using SHAP-style visualizations that explicitly link inputs to recommendations, improving comprehension and engagement when kept simple (Deo & Sontakke, 2021). Related literature highlights the risks of over-personalization and the role of transparent human–AI interaction (Capponi, Olafsson, & Zariphopoulou, 2022; Bianchi & Briere, 2021), while policy work emphasizes accountability and explainability in automated advice (Fein, 2017; Strzelczyk, 2017; European Commission, 2019). Overall, the literature posits transparency as a lever for trust and delegation, but it remains an empirical question whether profile-based disclosures change behavior in market settings.

2.3 Digital finance and the decision context

Digital finance reshapes how households decide, not just what they can access. By embedding decision support, reducing frictions, and shifting interactions onto smartphones, platforms alter attention, search, and default uptake. Empirically, broader digital access is associated with greater participation and risk-taking: exposure to digital services correlates with a higher allocation to risky assets and reduced risk aversion (Hong, Lu, & Pan, 2020); digital finance increases the extensive margin of risky asset holding (Shen, Hu, & Zhang, 2022). In China, a one-percent increase in digital finance (DFII) corresponds to a 0.13-percent rise in the share of risky assets and raises the likelihood of holding any risky assets (Hu, Guo, Shang, & Zhang, 2024). The interaction channel is therefore not neutral for portfolio choices.

Behaviorally, the mobile mindset favors fast, heuristic processing over deliberation (Lurie et al., 2018). In a lab setting with device randomization, smartphone users display a present-bias parameter about six percentage points lower and are roughly twelve points less willing to pay search costs for payoff-relevant information (Mograbi, 2022); habitual phone use similarly increases reliance on defaults and recommendations relative to larger screens (Wang, Malthouse, & Krishnamurthi, 2015). For explanation design, this implies that identical transparency features—such as profile-based disclosures of drivers and peer benchmarks can raise comprehension yet have device-contingent effects on acceptance: mobile may reinforce default adherence, whereas desktop affords

scrutiny and adjustment. This motivates our focus on desktop versus mobile sessions and frames our interpretation of heterogeneous impacts in the experiment.

3 Experimental design

3.1 Onboarding flow (“tunnel”)

In our setting, each prospective client who engages with the robo-advisor completes a standardized, multi-step onboarding process. For consistency with the partner’s internal terminology², we refer to this flow as the “tunnel” when discussing specific screens and steps. This structured sequence consists of four distinct stages, referred to as “Steps”, through which clients progress sequentially. Each step serves a specific function in the investment decision-making and account setup process (see Figure 1). In Step 1, a unique User ID is assigned and linked to the email account used to register, and the step concludes when the client initiates an investment simulation. In Step 2, the client completes a comprehensive financial questionnaire, which serves as the foundation for the risk assessment and portfolio recommendation process. Submission of this questionnaire simultaneously generates a new Contract ID linked to the corresponding User ID. Step 3 comprises the validation of the proposed risk profile, investment vehicle selection, and portfolio allocation against the client’s stated objectives and constraints. At Step 3, users may open a “Modify” panel listing alternative risk profiles with side-by-side allocations. Acceptance is recorded immediately after closing this panel when the final profile is confirmed and matches the originally recommended profile. Finally, Step 4 requires submission of all requisite legal documentation and the formal execution of the contractual agreement, thereby completing account opening with the robo-advisor.

A key component of this onboarding process is the financial questionnaire, which gathers comprehensive information regarding the client’s financial situation, investment preferences, and constraints. Specifically, the questionnaire includes inquiries on liquidity needs, risk tolerance, financial literacy, investment horizon, and investment objectives. Additionally, it collects various economic indicators such as financial wealth, income levels, savings capacity, and the amount allocated for the initial investment. Furthermore, socio-demographic variables are also recorded to provide a more holistic understanding of the investor profile. Upon completion of the questionnaire, the collected data are processed through a proprietary algorithm designed to generate a risk score. This score is instrumental in determining the client’s risk profile, which subsequently dictates the composition of their recommended investment portfolio. The algorithm not only assesses the appropriate level of risk exposure but also identifies the optimal investment vehicle. For the purposes of our study, we restrict attention to a tax-efficient saving plan (*assurance-vie*), the platform’s flagship product through which most clients invest. Focusing on a single investment vehicle ensures a common

²The industry partner’s identity is anonymized under an NDA that was a precondition for data access. Owing to confidentiality constraints and the partner’s deployment timeline, the study was not preregistered. The intervention and its implementation were co-designed with the partner to ensure product feasibility; however, the authors retained full independence over the empirical analysis and the write-up. The partner had no editorial control and no right to approve or veto publication.

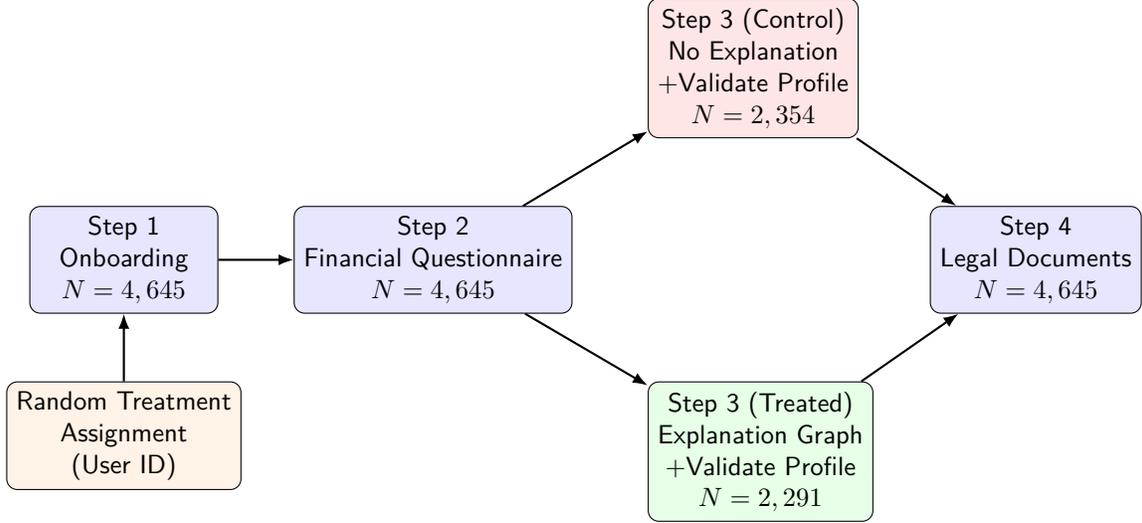


Figure 1: Flowchart of the Experiment

decision environment and maximizes sample size.

Table 1: Target asset allocation by risk profile (in %)

Asset Class	Risk Profiles								
	2	3	4	5	6	7	8	9	10
Euro Funds	70	60	40	20	0	0	0	0	0
Bond ETFs	15	20	30	40	50	40	30	20	0
Stock ETFs	15	20	30	40	50	60	70	80	100

Within the context of “*assurance-vie*” products, the generated risk score ranges from 2 to 10. A risk score of 2 corresponds to a highly conservative portfolio, allocated mostly to “Fonds Euro”(see Table 1). This investment vehicle provides capital protection and is characterized by low risk, albeit with limited return potential. At the opposite end of the risk spectrum, a score of 10 represents a fully equity-based portfolio, composed entirely of stocks, thus entailing a significantly higher level of risk. Intermediate risk profiles, with scores between 2 and 10, are defined by varying allocations between “Fonds Euro”, bonds, and equities. As the risk score increases, the proportion of “Fonds Euro” decreases, while the allocations to bonds and equities rise accordingly, reflecting a progressive shift towards a more aggressive investment strategy. Importantly, after receiving their assigned risk profile, prospective clients have the option to access a “modify option”, which allows them to compare their recommended profile against other available risk profiles. If desired, they may adjust their final risk selection before proceeding with investment validation and contract signing.

The intervention is an A/B test with randomization at Step 1 and treatment display at Step 3, immediately after completion of the Step 2 questionnaire. The treatment group receives an explanation of the assigned risk profile, which focuses solely on the rationale behind the score

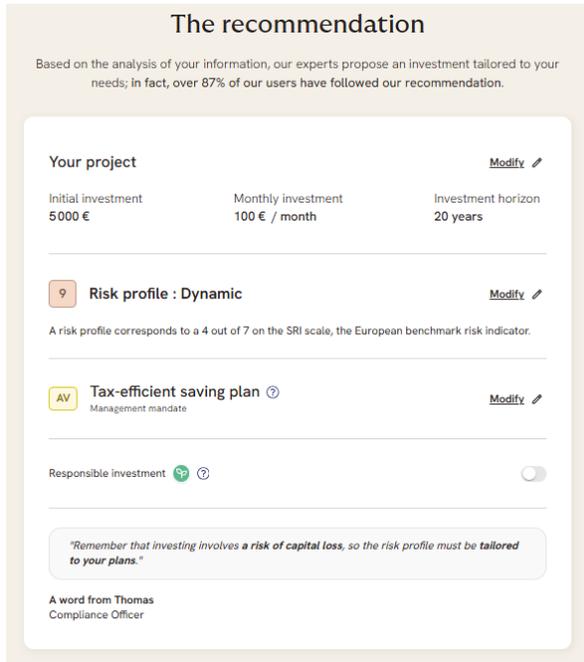
generated by the algorithm. In contrast, the baseline group does not receive any explanatory information. The treatment is randomly assigned at the beginning of Step 1 based on the User ID. Consequently, if a prospective client undergoes the onboarding process multiple times, a frequent occurrence in practice, they will consistently receive the same treatment condition. This design ensures that exposure to the intervention remains stable, thereby minimizing selection bias and preserving the integrity of the experimental framework. Randomization achieved balance across pre-treatment covariates; Mann–Whitney and χ^2 tests fail to reject equality across arms (all $p > 0.11$; Bonferroni-adjusted $p \geq 0.47$, see Table 5 in Appendix B).

3.2 Explanation Methodology

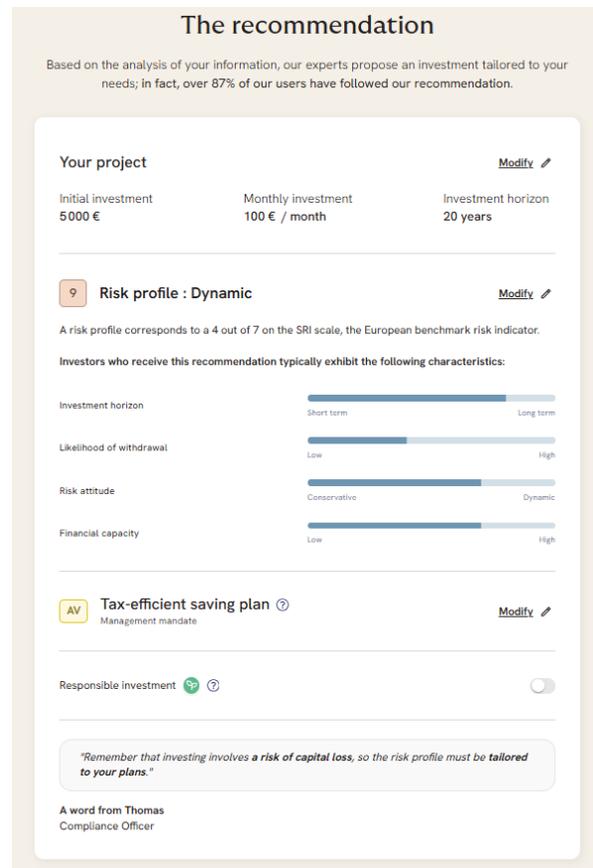
In this study, we focus on the transparency of the algorithm’s recommendations generated by a financial algorithm embedded within a robo-advisor. The algorithm computes a risk profile score for prospective clients, which in turn informs the investment recommendations they receive. To improve transparency and user understanding of the profile assessment process, we implement a graphical representation that visually communicates the relative importance of various factors contributing to the profile score. Specifically, the bar chart displays normalized weights for key determinants (liquidity constraints, risk preferences, economic characteristics, investment horizon) on a 0–10 scale. For visualization, we truncate displayed bars to the range 1–9 to avoid saturation at extremes while not altering relative rankings. To aid interpretation, the axis is annotated with attribute-specific verbal anchors (e.g., "short-term" and "long-term" for horizon; "low" and "high" for liquidity constraints; "conservative" and "dynamic" for risk attitude; "low" and "high" for financial capacity), which describe the underlying dimension rather than the importance units. An example of the visual explanation shown to treated users is provided in Figure 2b. The figure illustrates the characteristics of a user assigned risk profile 9. The accompanying bars show how individuals with similar profiles typically score on attributes such as investment horizon, withdrawal likelihood, risk attitude, and financial capacity.

The graphical content is tailored to each client’s assigned risk profile, which ranges from 2 to 10. Crucially, the explanation is profile-based: for each risk profile, the bar chart illustrates the mean normalized importance of each determinant across all users assigned that same score. In other words, it shows what characterizes the "typical" individual in that risk category, offering users comparative insights into what the algorithm considers the average user within this risk class. Displayed "importance" reflects model weights and is not a causal attribution. This approach embeds the explanation within a broader reference frame and seeks to enhance interpretability by contextualizing the recommendation with feature-importance cues linked to the recommended risk profile.

The financial algorithm underlying the robo-advisor aggregates responses from a financial questionnaire to compute a global risk score. Each response contributes to the final score based on a predefined weighting system that reflects the relative importance of various attributes in assessing a client’s investment profile. Although the scoring process is non-linear—such that small changes in



(a) Recommendation Interface for Non-Treated Users, translated from French



(b) Recommendation Interface for Treated Users, translated from French

inputs do not always produce proportional changes in the output—it is internally transparent: the structure and logic of the model are fully observable to the platform designers. In other words, the algorithm qualifies as a “white-box” model from a backend perspective. The display is a faithful transformation of the production scoring logic with no surrogate or post-hoc explainer used. This internal transparency ensures that the graphical representation used in the intervention remains an accurate reflection of the underlying risk assessment process, allowing for a reliable interpretation of how various factors influence the final investment profile. However, this interpretability is not directly available to end users, who only see the final score without insight into how their responses shaped it. The graphical explanation provided in our intervention serves to bridge this gap by offering users an interpretable summary of the main factors influencing their assigned risk profile, thereby reflecting the logic of the underlying model in an accessible and meaningful way to prospective clients.

Our approach builds on established techniques in the XAI literature, particularly feature importance attribution methods that highlight which inputs drive model outputs. While our algorithm is not a black-box model requiring post-hoc methods like SHAP, the graphical explanation serves a similar function: to make the recommendation logic more transparent and intuitive. By deliv-

ering profile-based explanations grounded in real user profiles, our design aims to foster greater comprehension and trust in algorithmic recommendations, contributing to the broader literature on interpretable decision support in finance.

3.3 Outcome variables

We exploit a rich array of outcome measures, which we organize into two broad categories: portfolio choices, capturing how users made decisions once they receive a risk-score recommendation, and engagement outcomes, reflecting the intensity and nature of interaction with the “Tunnel” interface.

Portfolio choices focus on (i) the acceptance of the recommended risk profile; (ii) the deviation, which retains the sign of the deviation to distinguish upward from downward adjustments; and (iii) separate positive and negative deviations to gauge the intensity of upward versus downward modifications.

Engagement outcomes are inspired by the XAI and recommender-systems literatures and quantify user interactions with the platform: (i) the total number of attempts across all steps (whether or not a contract was ultimately signed); (ii) the number of opening events, defined as instances in which users clicked to access the modification interface; and (iii) the number of modification events, defined as instances in which users actively adjusted their risk profile within that interface.

Further details on the distribution of these outcome variables can be found in Table 6 in Appendix B.

3.3.1 Portfolio Choices

Acceptance of score recommendation We define a binary variable $Accept_i$ that equals 1 if user i selects exactly the risk score recommended by the robo-advisor, and 0 otherwise. This outcome measures users’ willingness to follow the algorithmic suggestion.

Deviation measures from the recommended score We capture three related measures of how users adjust their chosen score relative to the robo-advisor’s recommendation.

First, the Deviation D_i preserves the sign of the difference—positive for upward adjustments and negative for downward adjustments—thereby indicating the deviation direction in risk-profile notches. Formally:

$$D_i = \text{ChosenRP}_i - \text{RecommendedRP}_i.$$

In this measure, when clients chose a score that is higher (lower) than the recommended score, the deviation is recorded using a positive (negative) value. This allows for documenting the direction of deviation from the recommended score, accounting for the sign of the deviation, if there is one.

Second and third, we subset signed deviations into separate measures of upward and downward

modifications, respectively denoted PD_i and ND_i . Then, on the subsample with $D_i > 0$ we define

$$PD_i = D_i \quad (\text{for } D_i > 0),$$

and, on the subsample with $D_i < 0$ we define

$$ND_i = D_i \quad (\text{for } D_i < 0).$$

3.3.2 Engagement Outcomes

Total number of attempts We define $TotAttempts_i$ as the total count of times user i completed the Tunnel until Step 3, irrespective of whether a contract was ultimately signed. This aggregate metric captures overall engagement with the onboarding process.

Number of opening events We define $OpenEvents_i$ as the number of instances in which user i clicked to access the modification interface. Each opening event reflects a user’s decision to review or reconsider the recommended risk profile.

Number of modification events We define $ModEvents_i$ as the number of instances in which user i actively adjusted their risk profile within the modification interface. This outcome measures the intensity of interaction with the recommendation system.

4 Results

4.1 Descriptive statistics

During the 105-day field experiment, 4,645 contracts³ were signed, corresponding to $N = 3,856$ unique users, with an average of 1.20 contracts per user (SD = 0.57). Of these, 2,354 belonged to the control group and 2,291 to the treatment group. The mean age of users was 38.00 years (SD = 13.00 years). Annual income was distributed as follows: 25.67% reported €30,000–45,000; 20.82% €60,000–100,000; 19.73% €45,000–60,000; 19.24 % < €30,000; 10.66% €100,000–150,000; and 3.86% > €150,000. Based on financial-knowledge items, 16.8% of users were classified as “*Non initié*” (uninitiated), 43.8% as “*Néophyte*” (beginner), and 39.5% as “*Sachant*” (knowledgeable). Finally, 86.87% of users were deemed novice in financial experience, while 13.12% were deemed experienced.

Baseline balance tests reveal no significant differences between the control and treatment groups in the distribution of recommended risk profiles (two-sided Mann–Whitney U test, p -value = 0.746; see Figure 3) nor in the distribution of device sessions (personal computer vs. phone; two-sided Mann–Whitney U test, p -value = 0.614; see Figure 10 in Appendix B). Additionally, there were

³For confidentiality reasons, the total number of users who entered the Tunnel and the exact start and end dates of the experiment cannot be disclosed. Consequently, our analysis focuses on those who signed a contract and is conducted at the contract level.

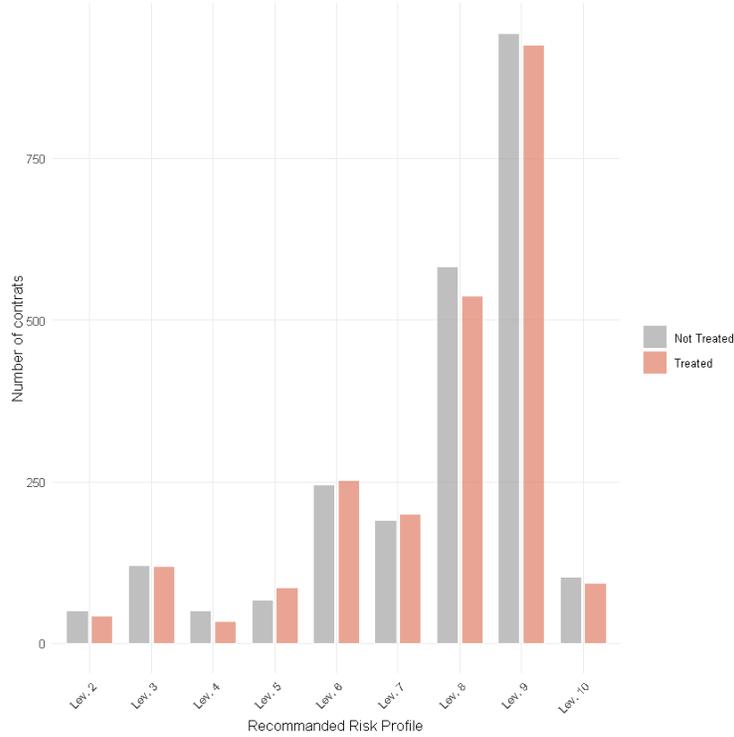


Figure 3: Distribution of Recommended Risk Profile by Experimental Condition

no differences in the distribution of user type — “first-time users”, who had not used the platform before the experiment, versus “returning users”, who the User ID was created beforehand and thus had used the platform previously (two-sided Mann–Whitney U test, p -value = 0.742; see Figure 11 in Appendix B).

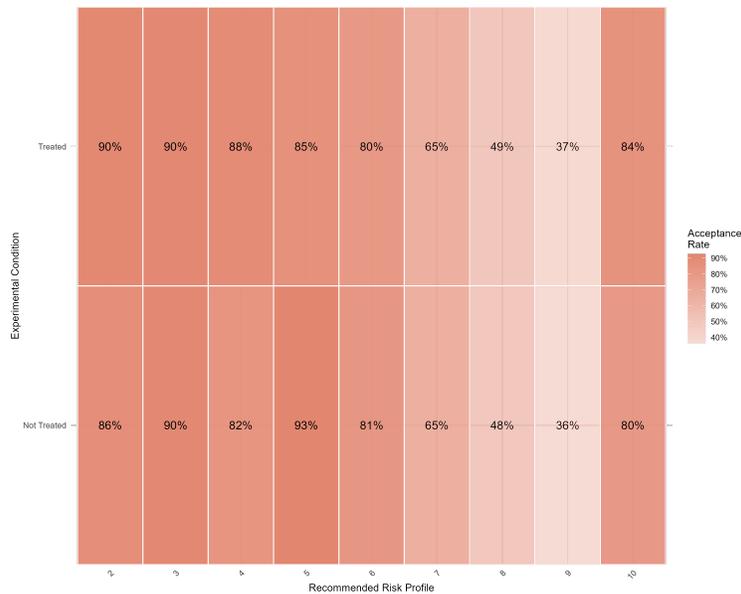
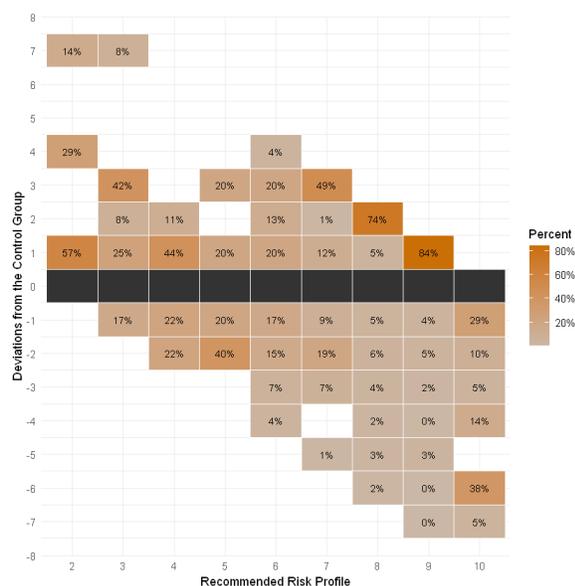
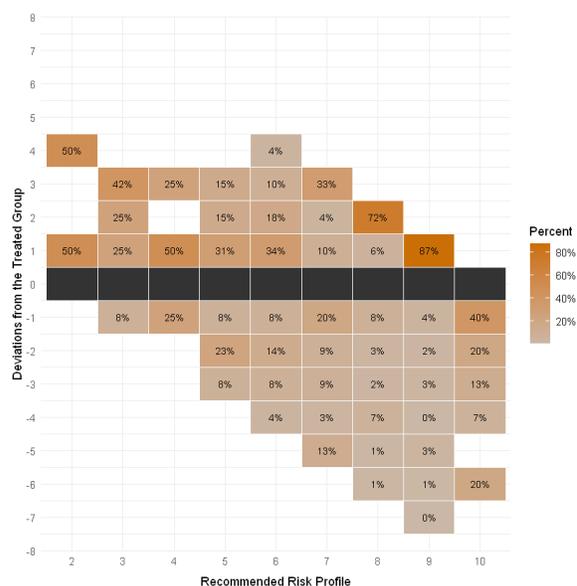


Figure 4: Acceptance Rate per Recommended Risk Profile and Experimental Condition

Descriptive statistics indicate that the overall acceptance rate does not differ significantly between the two experimental conditions: the control group exhibits a 54.63% acceptance rate, while the treatment group shows a 55.26% acceptance rate (two-sided Mann–Whitney U test, $p = 0.666$). Figure 4 displays the mean acceptance rate by risk profile for each group, revealing that both groups follow an almost identical pattern: acceptance rates range between approximately 80% and 90% for risk profiles 2 through 6 and profile 10, decline to below 70% for profile 7, drop below 50% for profile 8, and fall below 40% for profile 9. Notably, risk profiles 7–9 account for roughly 72% of our observations (see Table 7 in Appendix B), underscoring the practical importance of these lower acceptance segments.



(a) Deviations' Heatmap: Control Group



(b) Deviations' Heatmap: Treated Group

From the 2,093 observations in which users chose a different risk profile than recommended (1,645 positive deviations and 448 negative deviations), descriptive statistics reveal broadly similar deviation patterns across the two experimental conditions. The two heatmaps represent, for each recommended risk level, the percentage of users whose chosen profile lies at each deviation from the control group (Figure 5a) and the treated group (Figure 5b). In both graphs, deviations of +1 and +2 (i.e., selecting a risk profile one or two profiles higher than initially recommended) dominate at mid-range recommendations (profiles 3–7). For example, in the non-treated group, 42 percent of users recommended profile 3 actually selected profile 4 (+1 deviation), while 49 percent of users recommended profile 6 opted for profile 7 (+1). Similarly, in the treated group, 42 percent of profile 3 users and 33 percent of profile 6 users deviated upward by one notch. Upward deviations peak at recommended profile 8: 74 percent of non-treated users and 72 percent of treated users who deviated from it selected profile 10 (+2 deviation from 8). Notably, a substantial fraction of users recommended profiles 7 and 9 also deviated all the way to profile 10, suggesting they may have already had a predetermined preference for a fully equity-based portfolio. Negative deviations (selecting a lower risk profile than initially recommended) are relatively rare for mid-range

recommendations but become more prevalent at the highest risk levels. For instance, among users recommended profile 10, 38 percent of treated users and 29 percent of non-treated users who deviated chose profile 4 (−6 deviation from 10), reflecting a strong downward shift for the riskiest recommendation. These patterns mirror the summary statistics in Table 8: both groups exhibit very similar weighted mean deviations for several recommended profiles (e.g., Profile 6 shows means of 0.367 in the non-treated group versus 0.400 in the treated group, and Profile 8 shows 0.907 versus 0.876), with modes typically at +1 for mid-range profiles.

4.2 Treatment Effect

Table 2: Explanations Effect on the Acceptance Rate

	<i>Dependent variable:</i>			
	<i>Accept_i</i>			
	(1)	(2)	(3)	(4)
Treated	0.014 (0.073)	0.049 (0.106)	0.016 (0.073)	−0.068 (0.118)
Phone	0.314*** (0.078)	0.348*** (0.108)	0.317*** (0.078)	0.319*** (0.078)
Treated × Phone		−0.069 (0.148)		
First-Timers			0.327*** (0.119)	0.259* (0.141)
Treated × First-Timers				0.139 (0.151)
Constant	2.458** (1.080)	2.437** (1.081)	2.119* (1.090)	2.161** (1.090)
Controls	Yes	Yes	Yes	Yes
Observations	4,645	4,645	4,645	4,645
Log Likelihood	−2,567.506	−2,567.380	−2,562.848	−2,562.366
Pseudo R ² (McFadden)	0.197	0.197	0.198	0.198

Note: Clustered robust standard errors by user ID are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy and the initially recommended risk profile. . *p<0.1; **p<0.05; ***p<0.01

Table 2 reports the results of logistic regressions on the likelihood of users accepting the Robo-Advisor’s recommended risk profile (*Accept_i*), pooling all experimental conditions (Models 1–4). In each specification, we include a comprehensive set of socio-demographic and invest-

ment-related controls (as detailed in the note). Model 1 shows that the treatment indicator (i.e., provision of an explanation) has a small, positive coefficient ($\beta = 0.014$), but is not statistically significant ($SE = 0.073$). Across Models 2–4, we introduce interaction terms between the treatment indicator and two user characteristics (detailed in Section 4.1): device session (Phone) and user type (First-Timers). In Model 2, the interaction between Treatment and Phone is negative ($\beta = -0.069$) but statistically insignificant ($SE = 0.148$), leaving the marginal effect of explanations on phone users negligible. In Model 3, the main effect of explanations remains insignificant and in Model 4, which adds the Treatment \times First-Timers interaction, both treatment and the interaction terms are statistically indistinguishable from zero ($\beta = -0.068$, $SE = 0.118$; $\beta = 0.139$, $SE = 0.151$). Throughout all specifications, the Treatment coefficient does not attain significance, indicating that providing explanations does not meaningfully influence acceptance rates among users interacting with the same robo-advisor. By contrast, device session has a robust positive association with acceptance: users on mobile phones have approximately 37% higher odds to accept the recommendation than users on a computer (Model 1: $\beta = 0.314$, $SE = 0.078$, $p < 0.01$), an effect that persists across all specifications. Beyond device type, platform experience also matters. First-time users exhibit meaningfully higher acceptance than returning users (Model 3: $\beta = 0.327$, $SE = 0.119$, $p < 0.01$), corresponding to roughly 39% higher odds of accepting the recommendation; this attenuates to about 30% when the Treatment \times First-Timers interaction is included (Model 4: $\beta = 0.259$, $SE = 0.141$, $p < 0.10$). Taken together, explanations do not measurably change acceptance, while device type and platform experience are the primary correlates of adherence.

Table 3 presents OLS results for deviations (D_i) using again the 2,093 cases where users departed from the recommended profile. All four models control for demographics, investment characteristics, and initially recommended risk. Across all specifications, the treatment coefficient is small and not statistically different from zero. Likewise, being on a mobile device or a first-time user (and their interactions with treatment) shows no reliable impact on whether users shift toward higher or lower risk portfolios. The positive constant sign reflects the overall tendency to increase risk (1,645 positive vs. 448 negative deviations), but findings show that providing explanations does not shift that underlying direction.

Table 9 in Appendix B reports OLS estimates for positive deviation (PD_i), including only the 1,645 cases where users chose a higher risk profile than recommended. All four regressions control for the same socio-demographic, investment variables, and initially recommended risk. In Model 1, the explanation treatment has a small, negative coefficient ($\beta = -0.022$, $SE = 0.018$) that is not statistically significant, indicating no appreciable shift in upward deviations. Adding the Treatment \times Phone interaction in Model 2 does not change this result ($\beta = -0.024$, $SE = 0.027$; interaction $\beta = 0.003$, $SE = 0.040$). Likewise, Model 3 and 4 also shows that treated users exhibit a slightly smaller upward deviation (3, $\beta = -0.022$, $SE = 0.018$; 4, $\beta = -0.024$, $SE = 0.027$) but without significance, and Model 4’s Treatment \times First-Timers interaction is effectively zero ($\beta = 0.003$, $SE = 0.038$). Across all specifications, no coefficient on treatment, device, or user type meaningfully alters the magnitude of upward shifts.

Table 3: Explanations Effect on Deviations

	<i>Dependent variable:</i>			
	D_i			
	(1)	(2)	(3)	(4)
Treated	0.040 (0.078)	-0.029 (0.106)	0.039 (0.078)	0.093 (0.121)
Phone	-0.028 (0.079)	-0.098 (0.111)	-0.028 (0.079)	-0.030 (0.079)
Treated \times Phone		0.144 (0.156)		
First-Timers			-0.028 (0.117)	0.020 (0.138)
Treated \times First-Timers				-0.100 (0.157)
Constant	0.702 (1.099)	0.757 (1.107)	0.733 (1.102)	0.693 (1.098)
Controls	Yes	Yes	Yes	Yes
Observations	2,093	2,093	2,093	2,093
Adjusted R ²	0.221	0.221	0.221	0.221

Note: Clustered robust standard errors by user ID are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy and the initially recommended risk profile. * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

Table 4 focuses on the 448 instances where users moved to a lower-risk portfolio than recommended, estimating how explanations affect the magnitude of these downward shifts. In Model 1, the treatment coefficient is negative and non-significant (Model 1: $\beta = -0.244$, SE = 0.150). Adding the Treated \times Phone term in Model 2 increases both the magnitude and significance of the treatment coefficient (Model 2: $\beta = -0.452$, SE = 0.189, $p < 0.05$), but the interaction itself is not significant, suggesting that treatment effects vary across devices. While there is no statistical differences between phone users, it seems that treated computer users exhibit larger downward moves than their untreated counterparts. However, once we account for platform experience, the treatment effect loses significance (Model 3: $\beta = -0.137$, n.s.; Model 4: $\beta = 0.091$, n.s.), whereas first-time users exhibit markedly smaller downward deviations (≈ 0.79 – 0.98 , $p < 0.01$). Complementary robustness checks adding the Treated \times Phone \times First-Timers interaction yield no reliable differential

Table 4: Explanations Effect on Negative Deviations

	<i>Dependent variable:</i>			
	ND_i			
	(1)	(2)	(3)	(4)
Treated	-0.244 (0.150)	-0.438** (0.188)	-0.137 (0.148)	0.091 (0.241)
Phone	-0.059 (0.162)	-0.277 (0.223)	-0.049 (0.160)	-0.071 (0.159)
Treated \times Phone		0.449 (0.324)		
First-Timers			0.787*** (0.214)	0.983*** (0.266)
Treated \times First-Timers				-0.424 (0.304)
Constant	-1.488 (0.921)	-1.368 (0.906)	-3.098*** (1.051)	-3.207*** (1.052)
Controls	Yes	Yes	Yes	Yes
Observations	448	448	448	448
Adjusted R ²	0.134	0.136	0.156	0.158

Note: Clustered robust standard errors by user ID are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy and the initially recommended risk profile.. *p<0.1; **p<0.05; ***p<0.01

treatment effect by device or tenure (Appendix Table 10). We therefore interpret the desktop-specific amplification seen in simpler models as sensitive to specification; after accounting for user experience, there is no robust evidence that explanations increase the size of downward deviations.

In summary, transparency does not increase acceptance and does not affect deviations. The salient systematic differences are by device: desktop users accept less and deviate more than phone users. In simpler specifications, explanations are associated with larger downward moves among desktop users who already choose safer-than-recommended portfolios; however, this association becomes statistically insignificant once we account for user experience, and robustness checks with higher-order interactions do not reveal a stable differential effect by device or tenure. Taken together, the evidence points to behavior being shaped primarily by device context and platform experience, rather than by the explanation per se: adherence is higher on mobile and among first-time users; desktop users are more inclined to adjust, but explanations do not reliably push them further. Additionally, the overall deviation patterns, particularly the large shifts among those

moving to a full-equity portfolio, suggest that many users may have already held a predetermined investment preference, rendering the robo-advisor’s recommendation largely irrelevant regardless of treatment.

4.3 Additional analyses

Engagement outcomes Across all three engagement measures—total attempts (Table 11, Appendix B), “open” clicks to access the modification page (Table 12, Appendix B), and actual modification events (Table 13, Appendix B)—providing a profile-based graphical explanation had no impact on user behavior. In models predicting the total number of attempts, users who received an explanation neither reentered the tunnel process more nor fewer times than those in the control group, suggesting that explanations did not alter users’ persistence or hesitancy during the onboarding process. Similarly, when examining how often participants clicked on the option to open the modification page, the only hint of a treatment effect was a small increase among treated desktop users when compared to untreated desktop users; however, this increase does not remain statistically significant once mobile usage and other covariates were taken into account, indicating no robust difference in exploratory behavior attributable to explanations. Finally, the count of actual modification events remained unchanged, implying that explanations did not spur users to explore other profiles by revising their choices more frequently. Engagement nonetheless differs systematically by device: phone sessions exhibit fewer openings of the modification option and fewer recorded modifications than desktop sessions (Tables 12–13; baseline coefficients around -0.34 for opens and -0.37 for modifications). This device pattern aligns with our main behavioral results: mobile users are more likely to accept the initial recommendation with minimal exploration, whereas desktop users explore more. Moreover, platform experience is a second, strong correlate of engagement. First-time users are markedly more active across all three measures, with large positive coefficients for attempts (≈ 0.73 – 0.73), opens (≈ 0.84 – 0.87), and modifications (≈ 0.75 – 0.85), all significant at $p < 0.01$ (Tables 11–13). Collectively, these findings demonstrate that, contrary to expectations from human–computer interaction and XAI research, displaying how the robo-advisor positions each user relative to the “average” risk profile neither increased nor decreased engagement. In other words, revealing prototypical characteristics for each risk class did not lead participants to explore more options, revisit the onboarding process more often, or investigate alternate portfolios. This suggests that engagement on this platform is shaped by factors beyond algorithmic transparency, such as the type of device used and user experience.

Treatment heterogeneity Across both outcomes—acceptance and signed deviation—flexible CATE estimators confirm a near-zero average effect with only modest tails. Methodologically, we follow the framework in (Jacob, 2021): (i) Causal Random Forests grown with “honest” sample-splitting (one subsample to form leaves, a held-out subsample to estimate within-leaf treatment and control means), and (ii) Causal BART, which models the response surface as a sum of weak trees under regularizing priors and uses MCMC to obtain posterior draws of individual treatment

effects. In our randomized setting, identification is straightforward (unconfoundedness holds by design); we nonetheless use honest splitting/held-out estimation to limit adaptive overfitting and obtain uncertainty via out-of-bag variance estimates for the forest and posterior credible intervals for BART. Tuning (number/depth of trees, priors, burn-in) follows (Jacob, 2021).

For acceptance, both methods yield tightly centered CATE distributions with ATEs around 0 with small positive and negative fringes (Figures 12a-12b; Table 14 in Appendix B). For deviation, both estimators produce CATE distributions that are still tightly centered near zero, with only modestly wider tails than for acceptance (Figures 13a-13b; Table 16 in Appendix B). The only consistent moderator is project objective: savings goals align with slightly more upward moves, retirement goals with slightly more downward moves (Table 17, Appendix B).

5 Conclusion

This paper uses a randomized controlled trial embedded in the interface of a leading French robo-advisor to test whether profile-based graphical explanations influence portfolio choices. Two facts stand out. First, contrary to prevailing expectations in the XAI and behavioral finance literatures, explanations neither increase adherence to the recommended risk score nor reduce the magnitude of deviations when they opt to choose a different profile. Additionally, engagement metrics, such as the number of tunnel completions, interface openings, or actual profile modifications, remain statistically indistinguishable across treatment and control groups. Second, behavior varies systematically by context: relative to phone sessions, desktop sessions accept less and explore more, and first-time users accept more and engage more than returning users. In some simpler specifications we observe larger downward moves among desktop users who already choose a safer-than-recommended portfolio; however, this pattern disappears and becomes statistically indistinguishable from zero once we account for user experience, and additional robustness checks do not reveal a stable differential treatment effect by device or tenure.

These null results carry important implications for both theory and practice. First, they suggest that in high-stakes financial decisions, where many investors appear to hold ex-ante portfolio preferences, algorithmic explanations may be insufficient to overcome deeply held convictions or “algorithm aversion.” Second, they highlight the need to reconsider the design of explainability interventions: peer-based benchmarks alone may not resonate with end-users unless coupled with interactive, personalized, or narrative components that address individual concerns more directly. Third, from a policy standpoint, our findings caution regulators and platform designers against assuming that transparency mandates will automatically translate into greater adoption or trust without complementary measures.

If explanations matter, they appear to do so through context rather than on average. Desktop interactions are associated with more deliberation and adjustment; mobile interactions with higher baseline adherence and fewer adjustments. But once user experience is taken into account, we do not find robust evidence that explanations systematically push desktop users further toward lower

risk. The design takeaway is therefore not “add explanations everywhere,” but “tailor presentation to context and experience.”

For design and policy, the implication is to favor context-contingent transparency over one-size-fits-all. Platforms aiming to preserve adherence at the acceptance step could test lighter or deferred explanations on desktop (e.g., progressive disclosure, post-confirmation details, or optional “learn more”), evaluate alternative benchmark framings, and segment by tenure. Conversely, mobile flows may not warrant additional surface area for explanations at the decision point. More broadly, transparency alone may shift adjustments rather than the decision to accept itself.

Our setting bounds inference. We study a white-box, peer-based explanation during onboarding for life-insurance wrappers, observe short-run choices (not funding, rebalancing, or performance), and evaluate one market and design. Other explanation styles (e.g., counterfactuals, conversational agents) or products may behave differently, and longer horizons may reveal delayed trust effects.

Overall, the evidence maps boundary conditions for XAI in digital financial advice: no average improvement in acceptance, behavior shaped primarily by device context and platform experience, and only weak, specification-sensitive signs that explanations could magnify downward moves on desktop. While algorithmic transparency remains an essential component of accountable FinTech, effective explainability will likely come from adaptive, device and experience-focused designs when adherence is desired; richer and interactive when exploration is the goal rather than uniform displays of feature importance. Future work should test such adaptive frameworks across investor segments and over extended horizons.

Acknowledgments

This research has been conducted within the PREF research initiative under the aegis of the Europlace Institut of Finance, a joint initiative by YOMONI. It was also supported by the APR-IA AcceptAlgo, funded by the Centre-Val de Loire Region. Finally, it was supported by the ANR AESOP.

References

- Abraham, F., Schmukler, S. L., & Tessada, J. (2019, Feb). *Robo-advisors : Investing through machines* (Research and Policy Briefs No. 134881). The World Bank.
- Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6, 52138–52160.
- Anjomshoae, S., Najjar, A., Calvaresi, D., & Främling, K. (2019). Explainable agents and robots: Results from a systematic literature review. In *18th international conference on autonomous agents and multiagent systems (aamas 2019), montreal, canada, may 13–17, 2019* (pp. 1078–1088).

- Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... others (2020). Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion*, 58, 82–115.
- Beaudouin, V., Bloch, I., Bounie, D., Cléménçon, S., d’Alché Buc, F., Eagan, J., ... Parekh, J. (2020). Flexible and context-specific ai explainability: a multidisciplinary approach. *arXiv preprint arXiv:2003.07703*.
- Ben David, D., Resheff, Y. S., & Tron, T. (2021). Explainable ai and adoption of financial algorithmic advisors: an experimental study. In *Proceedings of the 2021 aaai/acm conference on ai, ethics, and society* (pp. 390–400).
- Bianchi, M., & Brière, M. (2020). Robo-advising for small investors. *SSRN Electronic Journal*.
- Bianchi, M., & Briere, M. (2021). Robo-advising: less ai and more xai? *Available at SSRN 3825110*.
- Cai, C. J., Jongejan, J., & Holbrook, J. (2019). The effects of example-based explanations in a machine learning interface. In *Proceedings of the 24th international conference on intelligent user interfaces* (pp. 258–262).
- Capponi, A., Olafsson, S., & Zariphopoulou, T. (2022). Personalized robo-advising: Enhancing investment through client interaction. *Management Science*, 68(4), 2485–2512.
- Cardillo, G., & Chiappini, H. (2024). Robo-advisors: A systematic literature review. *Finance Research Letters*, 62, 105119. Retrieved from <https://www.sciencedirect.com/science/article/pii/S1544612324001491> doi: <https://doi.org/10.1016/j.frl.2024.105119>
- Castelvecchi, D. (2016). Can we open the black box of ai? *Nature*, 538(7623), 20–23.
- Cheng, H.-F., Wang, R., Zhang, Z., O’connell, F., Gray, T., Harper, F. M., & Zhu, H. (2019). Explaining decision-making algorithms through ui: Strategies to help non-expert stakeholders. In *Proceedings of the 2019 chi conference on human factors in computing systems* (pp. 1–12).
- Cohen, S., Dror, G., & Ruppin, E. (2007). Feature selection via coalitional game theory. *Neural computation*, 19(7), 1939–1961.
- Covert, I., Lundberg, S., & Lee, S.-I. (2021). Explaining by removing: A unified framework for model explanation. *Journal of Machine Learning Research*, 22(209), 1–90.
- Cramer, H., Evers, V., Ramlal, S., Van Someren, M., Rutledge, L., Stash, N., ... Wielinga, B. (2008). The effects of transparency on trust in and acceptance of a content-based art recommender. *User Modeling and User-adapted interaction*, 18, 455–496.
- D’Acunto, F., & Rossi, A. G. (2023). Robo-advice: Transforming households into rational economic agents. *Annual Review of Financial Economics*, 15(1), 543–563.
- Deo, S., & Sontakke, N. S. (2021). Usability, user comprehension, and perceptions of explanations for complex decision support systems in finance: A robo-advisory use case. *Computer*, 54(10), 38–48.
- Dietvorst, B. J., Simmons, J. P., & Massey, C. (2015). Algorithm aversion: people erroneously avoid algorithms after seeing them err. *Journal of experimental psychology: General*, 144(1), 114.

- Dietvorst, B. J., Simmons, J. P., & Massey, C. (2018). Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them. *Management science*, *64*(3), 1155–1170.
- Eslami, M., Krishna Kumaran, S. R., Sandvig, C., & Karahalios, K. (2018). Communicating algorithmic process in online behavioral advertising. In *Proceedings of the 2018 chi conference on human factors in computing systems* (pp. 1–13).
- European Commission. (2019). *Ethics guidelines for trustworthy AI*. Retrieved from <https://ec.europa.eu/futurium/en/ai-alliance-consultation> (Available at: <https://ec.europa.eu/futurium/en/ai-alliance-consultation>)
- Fein, M. L. (2017). Are robo-advisors fiduciaries? *Available at SSRN 3028268*.
- Herlocker, J. L., Konstan, J. A., & Riedl, J. (2000). Explaining collaborative filtering recommendations. In *Proceedings of the 2000 acm conference on computer supported cooperative work* (pp. 241–250).
- Hong, C. Y., Lu, X., & Pan, J. (2020). *Fintech adoption and household risk-taking: From digital payments to platform investments* (Tech. Rep.). National Bureau of Economic Research.
- Hu, D., Guo, F., Shang, J., & Zhang, X. (2024). Does digital finance increase household risk-taking? evidence from china. *International Review of Economics & Finance*, *93*, 1197–1210.
- Jacob, D. (2021). Cate meets ml: Conditional average treatment effect and machine learning. *Digital Finance*, *3*(2), 99–148.
- Jung, D., Dorner, V., Weinhardt, C., & Pusmaz, H. (2018). Designing a robo-advisor for risk-averse, low-budget consumers. *Electronic Markets*, *28*(3), 367–380. doi: 10.1007/s12525-017-0279-9
- Krzyżniński, M., Spytek, M., Baniecki, H., & Biecek, P. (2023). Survshap (t): time-dependent explanations of machine learning survival models. *Knowledge-Based Systems*, *262*, 110234.
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in neural information processing systems*, *30*.
- Lurie, N. H., Berger, J., Chen, Z., Li, B., Liu, H., Mason, C. H., ... others (2018). Everywhere and at all times: mobility, consumer decision-making, and choice. *Customer Needs and Solutions*, *5*, 15–27.
- Mograbi, E. (2022). Decision-makers are more impulsive on smartphones than on computers. *Journal of Behavioral and Experimental Economics*, *100*, 101916.
- Morana, S., Gnewuch, U., Jung, D., & Granig, C. (2020). The effect of anthropomorphism on investment decision-making with robo-advisor chatbots. In *Ecis*.
- Nunes, I., & Jannach, D. (2017). A systematic review and taxonomy of explanations in decision support and recommender systems. *User Modeling and User-Adapted Interaction*, *27*, 393–444.
- Patel, K., & Lincoln, M. (2019). It’s not magic: Weighing the risks of ai in financial services. *Report available at <https://www.european-microfinance.org/publication/its-not-magic-weighing-risks-ai-financial-services>*.
- Philippon, T. (2019, 9). On fintech and financial inclusion. *National Bureau of Economic Research*.

doi: 10.3386/w26330

- Preece, A., Harborne, D., Braines, D., Tomsett, R., & Chakraborty, S. (2018). Stakeholders in explainable ai. *arXiv preprint arXiv:1810.00184*.
- Rader, E., Cotter, K., & Cho, J. (2018). Explanations as mechanisms for supporting algorithmic transparency. In *Proceedings of the 2018 chi conference on human factors in computing systems* (pp. 1–13).
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Why should i trust you?: Explaining the predictions of any classifier. *CoRR, abs/1602.04938*. Retrieved from <http://arxiv.org/abs/1602.04938>
- Riedl, M. O. (2019). Human-centered artificial intelligence and machine learning. *Human behavior and emerging technologies, 1*(1), 33–36.
- Shen, Y., Hu, W., & Zhang, Y. (2022). Digital finance, household income and household risky financial asset investment. *Procedia Computer Science, 202*, 244–251.
- Shin, D. (2021). The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable ai. *International Journal of Human-Computer Studies, 146*, 102551.
- Speith, T. (2022). A review of taxonomies of explainable artificial intelligence (xai) methods. In *Proceedings of the 2022 acm conference on fairness, accountability, and transparency* (pp. 2239–2250).
- Strzelczyk, B. E. (2017). Rise of the machines: The legal implications for investor protection with the rise of robo-advisors. *DePaul Bus. & Comm. LJ, 16*, 54.
- Wang, R. J.-H., Malthouse, E. C., & Krishnamurthi, L. (2015). On the go: How mobile shopping affects customer purchase behavior. *Journal of retailing, 91*(2), 217–234.
- Weitz, K., Schiller, D., Schlagowski, R., Huber, T., & André, E. (2021). Let me explain!: exploring the potential of virtual agents in explainable ai interaction design. *Journal on Multimodal User Interfaces, 15*(2), 87–98.

Appendix

A Experimental Design Appendices

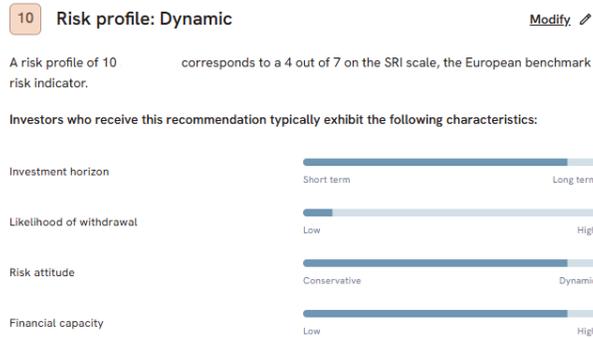


Figure 6: Graphical explanation for the Risk Profile 10



Figure 7: Graphical explanation for the Risk Profile 6

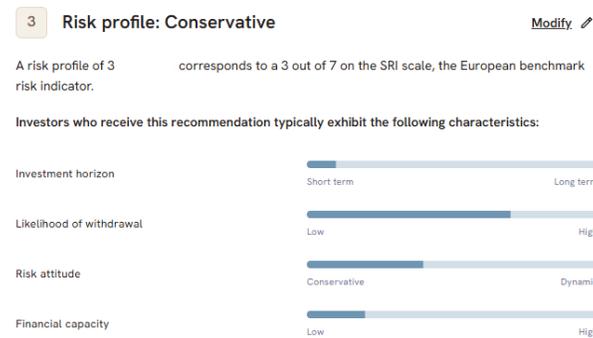


Figure 8: Graphical explanation for the Risk Profile 3

Modification of your allocation x

Your allocation can be modified at any time from your investor portal.

Responsible Investment ✔ 🔍 🔇

Selected risk profile: Profile 9

- Suravenir Opportunités 2 Euro Fund
- HSBC Clic Euro 85
- SC Real Estate
- Transparency III
- 2 Profile 2
- 3 Profile 3
- 4 Profile 4
- 5 Profile 5
- 6 Profile 6
- 7 Profile 7
- 8 Profile 8
- 9 Profile 9 Recommended for you
- 10 Profile 10

You want to take on significant risk with the prospect of high returns and you have a long-term investment horizon (>8 years). Profile 9 is therefore perfectly suited to your project.

Confirm

Figure 9: Modification Page

B Additional Results

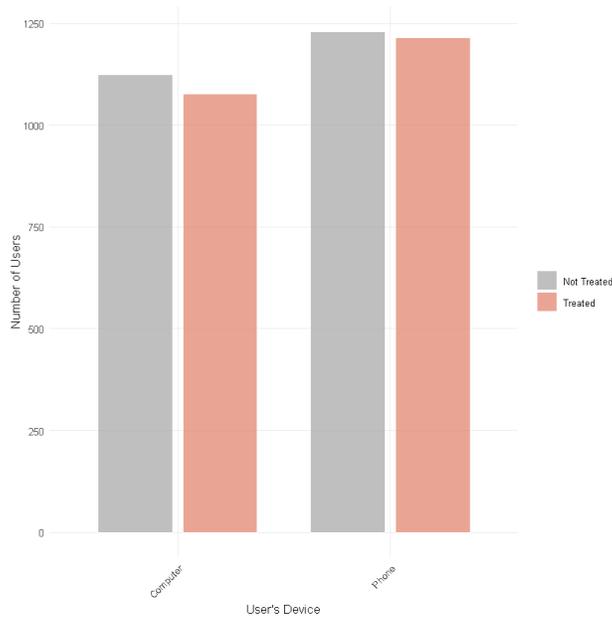


Figure 10: Distribution of Devices Used by Experimental Condition

Table 5: Balance tests between treatment and control groups (non-parametric)

Variable	<i>p</i> -value	Bonferroni <i>p</i>
<i>Numeric / ordinal covariates (Mann–Whitney)</i>		
User age	0.119	0.476
Initially declared investment amount	0.600	1.000
Rank of tunnel run	0.853	0.968
Initially recommended risk profile	0.749	1.000
<i>Categorical covariates (χ^2 with Fisher fallback)</i>		
Financial experience	0.625	1.000
Estimated Revenue	0.767	1.000
Type of investment project	0.571	1.000
Home-ownership status	0.872	1.000
Financial literacy	0.978	1.000
Investment Horizon	0.761	1.000

Notes: Two-sided tests. Numeric/ordinal covariates tested with Mann–Whitney (Wilcoxon rank-sum). Categorical covariates tested with Pearson χ^2 ; Fisher’s exact used when expected counts are small. Bonferroni adjustment applied within each family (numeric vs. categorical). Values rounded to three decimals.

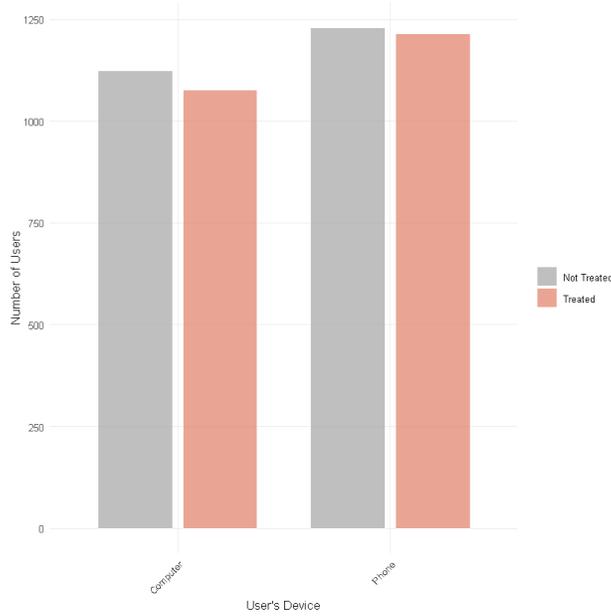


Figure 11: Distribution of Users' Type by Experimental Condition

Table 6: Summary Statistics for Behavioral and Engagement Outcomes

Outcome	Q1	Mean	Median	Q3	SD	N
Behavioral Outcomes						
Acceptance ($Accept_i$)	0.00	0.55	1.00	1.00	0.49	4,645
Deviation (D_i)	1.00	0.52	1.00	2.00	1.91	2,093
Positive deviation (PD_i)	1.00	1.39	1.00	2.00	0.64	1,645
Negative deviation (ND_i)	-4.00	-2.65	-2.00	-1.00	0.19	448
Engagement Outcomes						
Total attempts ($TotAttempts_i$)	1.00	1.48	1.00	2.00	0.99	3,856
Opening events ($OpenEvents_i$)	0.00	1.14	0.00	1.00	3.69	3,856
Modification events ($ModEvents_i$)	0.00	0.71	0.00	0.00	2.83	3,856

Table 7: Number of observations per recommended risk level by experimental condition and in total.

Risk Profile	Not Treated	Treated	Total	Overall (%)
Profile 2	51	42	93	2.0
Profile 3	121	119	240	5.2
Profile 4	51	34	85	1.8
Profile 5	67	87	154	3.3
Profile 6	245	252	497	10.7
Profile 7	191	200	391	8.4
Profile 8	582	538	1120	24.1
Profile 9	943	926	1869	40.2
Profile 10	103	93	196	4.2

Table 8: Summary Statistics of Deviation by Recommended Risk Level and Experimental Condition

Risk Profile	Non-treated			Treated		
	Total n	Weighted Mean	Mode	Total n	Weighted Mean	Mode
Profile 2	7	2.7143	1	4	2.5000	1
Profile 3	12	2.0833	3	12	1.9167	3
Profile 4	9	0.0000	1	4	1.0000	1
Profile 5	5	-0.2000	-2	13	0.3077	1
Profile 6	46	0.3696	1	50	0.4000	1
Profile 7	67	0.8507	3	70	-0.2143	3
Profile 8	300	0.9067	2	275	0.8764	2
Profile 9	601	0.4093	1	582	0.5206	1
Profile 10	21	-3.8095	-6	15	-2.6667	-1

Table 9: Explanations Effect on Positive Deviations

	<i>Dependent variable:</i>			
	<i>PD_i</i>			
	(1)	(2)	(3)	(4)
Treated	-0.022 (0.018)	-0.024 (0.027)	-0.022 (0.018)	-0.024 (0.027)
Phone	0.007 (0.019)	0.005 (0.031)	0.007 (0.019)	0.007 (0.019)
Treated × Phone		0.003 (0.040)		
First-Timers			-0.007 (0.027)	-0.008 (0.035)
Treated × First-Timers				0.003 (0.038)
Constant	2.136* (1.174)	2.137* (1.176)	2.143* (1.175)	2.144* (1.171)
Controls	Yes	Yes	Yes	Yes
Observations	1,645	1,645	1,645	1,645
Adjusted R ²	0.668	0.668	0.668	0.668

Note: Clustered robust standard errors by user ID are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy and the initially recommended risk profile. *p<0.1; **p<0.05; ***p<0.01

Table 10: Explanations Effect on Negative Deviations (Cont.)

	<i>Dependent variable:</i>		
	<i>ND_i</i>		
	(1)	(2)	(3)
Treated	-0.123 (0.270)	-0.131 (0.268)	-0.265 (0.317)
Phone	-0.302 (0.218)	-0.158 (0.288)	-0.295 (0.345)
Treated × Phone	0.480 (0.319)	0.471 (0.320)	0.751 (0.493)
First-Timers	0.990*** (0.263)	1.096*** (0.289)	0.988*** (0.314)
Treated × First-Timers	-0.406 (0.303)	-0.406 (0.303)	-0.176 (0.399)
Phone × First-Timers		-0.251 (0.288)	-0.004 (0.425)
Treated × Phone x First-Timers			-0.510 (0.595)
Constant	-3.105*** (1.040)	-3.260*** (1.053)	-3.176*** (1.066)
Controls	Yes	Yes	Yes
Observations	448	448	448
Adjusted R ²	0.160	0.159	0.158

Note: Clustered robust standard errors by user ID are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy and the initially recommended risk profile. *p<0.1; **p<0.05; ***p<0.01

Table 11: Explanations Effect on the Total Number of Attempts

	<i>Dependent variable:</i>			
	<i>TotAttempts_i</i>			
	(1)	(2)	(3)	(4)
Treated	-0.007 (0.022)	0.009 (0.030)	-0.001 (0.020)	0.002 (0.034)
Phone	0.012 (0.023)	0.027 (0.031)	0.013 (0.021)	0.013 (0.021)
Treated × Phone		-0.029 (0.042)		
First-Timers			0.726*** (0.031)	0.728*** (0.038)
Treated × First-Timers				-0.004 (0.041)
Constant	0.149 (0.143)	0.141 (0.145)	-0.577*** (0.137)	-0.578*** (0.137)
Controls	Yes	Yes	Yes	Yes
Observations	3,856	3,856	3,856	3,856
Log Likelihood	-5,042.480	-5,042.332	-4,912.462	-4,912.459
Pseudo-R ² (McFadden)	0.030	0.030	0.055	0.055

Note: Robust standard errors are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy, and the initially recommended risk profile. *p<0.1; **p<0.05; ***p<0.01

Table 12: Explanations Effect on the Total Number of Open Events

	<i>Dependent variable:</i>			
	<i>OpenEvents_i</i>			
	(1)	(2)	(3)	(4)
Treated	0.065 (0.101)	0.212* (0.123)	0.071 (0.101)	0.102 (0.171)
Phone	-0.340*** (0.094)	-0.171 (0.139)	-0.337*** (0.093)	-0.338*** (0.092)
Treated × Phone		-0.336* (0.195)		
First-Timers			0.839*** (0.139)	0.867*** (0.177)
Treated × First-Timers				-0.056 (0.209)
Constant	-1.084* (0.609)	-1.170* (0.609)	-1.970*** (0.658)	-1.989*** (0.664)
Controls	Yes	Yes	Yes	Yes
Observations	3,856	3,856	3,856	3,856
Log Likelihood	-7,918.582	-7,904.173	-7,779.042	-7,778.649
Pseudo-R ²	0.141	0.142	0.156	0.156

Note: Robust standard errors are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy, and the initially recommended risk profile. *p<0.1; **p<0.05; ***p<0.01

Table 13: Explanations Effect on the Total Number of Modification Events

	<i>Dependent variable:</i>			
	<i>ModEvents_i</i>			
	(1)	(2)	(3)	(4)
Treated	0.051 (0.117)	0.202 (0.142)	0.056 (0.117)	0.171 (0.189)
Phone	-0.366*** (0.108)	-0.191 (0.163)	-0.362*** (0.108)	-0.366*** (0.107)
Treated × Phone		-0.349 (0.224)		
First-Timers			0.747*** (0.174)	0.850*** (0.215)
Treated × First-Timers				-0.211 (0.234)
Constant	-1.689** (0.680)	-1.777*** (0.683)	-2.484*** (0.728)	-2.550*** (0.734)
Controls	Yes	Yes	Yes	Yes
Observations	3,856	3,856	3,856	3,856
Log Likelihood	-5,676.116	-5,666.519	-5,607.795	-5,604.267
Pseudo-R ²	0.152	0.154	0.162	0.163

Note: Robust standard errors are reported in parentheses. The regressions control for user age, type of investment project, investment horizon, initially declared investment amount, rank of tunnel run, estimated revenue, home-ownership status, financial experience, financial literacy, and the initially recommended risk profile. *p<0.1; **p<0.05; ***p<0.01

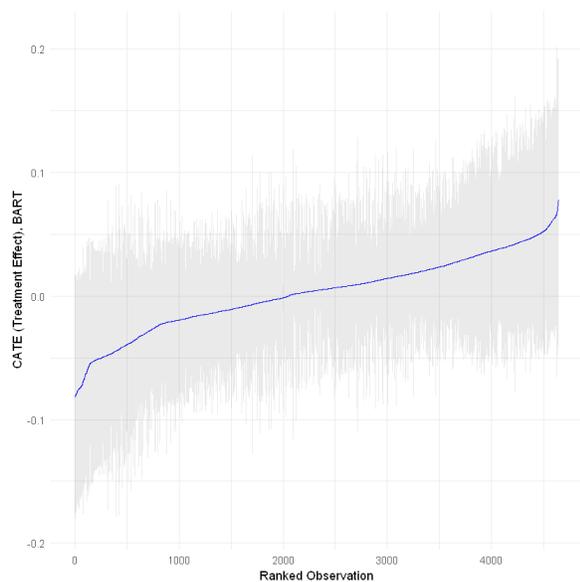
Table 14: CATE results for the Acceptance of the Recommendation

Method	20% least	ATE	20% most
Causal BART	-0.0381	0.0023	0.0429
Causal Forest	-0.0368	0.0059	0.0532

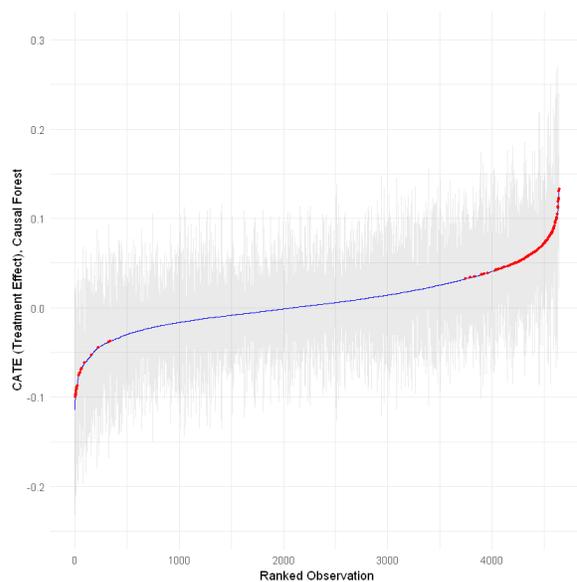
Table 15: Characteristics of Users by CATE Group, Acceptance of the Recommendation

Variable	No Significance (n=4,502)	Positive (n=119)	Negative (n=24)	p-value
<i>Proportions</i>				
Low patrimony (<30 k)	0.1966	0.0924	0.2917	0.0087**
Project Type: Savings	0.8827	0.9664	0.9167	0.0164*

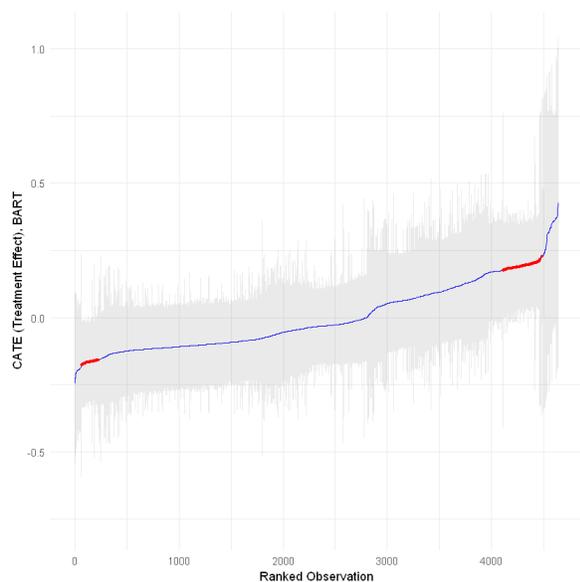
Note: Proportions are means of 0/1 indicators. p-values from one-way ANOVA across the three groups. *p<0.05; **p<0.01.



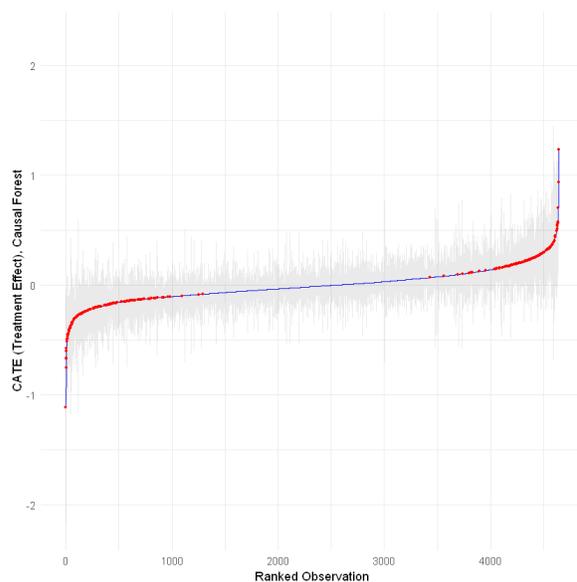
(a) Causal BART, Acceptance of the Recommendation



(b) Causal Forest, Acceptance of the Recommendation



(a) Causal BART, Deviation



(b) Causal Forest, Deviation

Table 16: CATE results for the Deviation

Method	20% least	ATE	20% most
Causal BART	-0.1368	0.0003	0.1979
Causal Forest	-0.1938	-0.005	0.2052

Table 17: User Characteristics by CATE Group and Method (Deviation)

Variable	BART				Causal Forest			
	NoSig	Pos	Neg	p(BART)	NoSig	Pos	Neg	p(CF)
<i>n</i>	4,347	210	88	—	4,454	53	138	—
<i>Proportions</i>								
Financial Knowledge: Beginner	0.445	0.409	0.329	0.063*	—	—	—	—
Financial Knowledge: Knowledgeable	0.393	0.433	0.500	0.068*	—	—	—	—
Investment Horizon	10.806	11.685	11.579	0.065*	—	—	—	—
Project type: Savings	0.883	0.933	0.863	0.068***	0.885	0.940	0.821	< 0.01***
Project type: Retirement	—	—	—	—	0.044	0.006	0.076	< 0.01***
<i>Means</i>								
User age (years)	—	—	—	—	21.54	17.90	23.22	< 0.01***
Recommended risk profile	6.615	6.985	6.909	< 0.01***	—	—	—	—

Note: “NoSig,” “Pos,” and “Neg” denote non-significant, positive, and negative CATE estimates, respectively; proportions for binary indicators, means for continuous variables. p-values from one-way ANOVA across the three CATE groups, by method. *p<0.10; **p<0.05; ***p<0.01.