

# Sell Me the Blackbox! Why eXplainable Artificial Intelligence (XAI) May Hurt Customers

Behnam Mohammadi<sup>1</sup>      Nikhil Malik<sup>2</sup>      Tim Derdenger<sup>3</sup>      Kannan Srinivasan<sup>4</sup>  
behnamm@cmu.edu   maliknik@usc.edu   derdenge@cmu.edu   kannans@cmu.edu

<sup>1,3,4</sup> Tepper School of Business, Carnegie Mellon University

<sup>2</sup> Marshall School of Business, University of Southern California

## Abstract

Recent AI algorithms are blackbox models whose decisions are difficult to interpret. eXplainable AI (XAI) seeks to address lack of AI interpretability and trust by explaining to customers their AI decision, e.g., decision to reject a loan application. The common wisdom is that regulating AI by mandating fully transparent XAI leads to greater social welfare. This paper challenges this notion through a game theoretic model for a policy-maker who maximizes social welfare, firms in a duopoly competition that maximize profits, and heterogeneous consumers. The results show that XAI regulation may be redundant. In fact, mandating fully transparent XAI may make firms and customers worse off. This reveals a trade-off between maximizing welfare and receiving explainable AI outputs. We also discuss managerial implications for policy-maker and firms.

**Keywords:** Machine Learning, Explainable AI, Economics of AI, Regulation, Fairness

## 1. Introduction

Recent years have seen a surge in the adoption of Artificial Intelligence (AI) models for decision-making. Gartner identifies AI engineering among the top 12 strategic technology trends of 2022<sup>1</sup>, and International Data Corporation (IDC) forecasts global spending on AI systems will jump from \$85.3 billion in 2021 to more than \$204 billion in 2025<sup>2</sup>. But a key challenge in adoption of AI is the interpretability of its decisions or predictions. While early AI models were easily interpretable, the latest methods such as Deep Neural Networks (DNNs) are opaque decision systems with very huge parametric spaces that make their decisions difficult to understand even by their developers<sup>3</sup>. In this regard, most recent AI algorithms are complex *blackbox* models (Castelvecchi, 2016). Sometimes, this is not a problem because the consequences may be negligible (e.g., email categorization AI), or because users may not want to know explanations for AI outputs (e.g., answers given by a voice assistant like Siri and

Alexa). But in most settings, humans are generally reluctant to adopt algorithms that are not interpretable, tractable, and trustworthy (Zhu et al., 2018).

To address AI interpretability, researchers have recently shifted their focus to eXplainable Artificial Intelligence (XAI), a class of methods that aim to produce “glass box” models that are explainable to humans while maintaining a high level of prediction accuracy<sup>4</sup> (Abdollahi and Nasraoui, 2016; Csiszár et al., 2020; Doshi-Velez and Kim, 2017; Holzinger et al., 2017b; Lipton, 2017; Murdoch et al., 2019). Put differently, the goal of XAI is to enable human users—including non-technical non-experts—to understand, trust, and effectively manage the emerging AI systems. XAI has been gaining traction across healthcare, retail, media and entertainment, and aerospace and defense. According to Explainable AI Market Report<sup>a</sup>, the market size of XAI across the globe was estimated to be \$4.4 billion in 2021 and is predicted to reach \$21.0 billion by 2030 with a compound annual growth rate (CAGR) of 18.4% from 2022 to 2030. But despite the growing interest in XAI—from researchers to engineers and governments<sup>5</sup>—little is known about the *economic implications* of XAI for firms and consumers (Adadi and Berrada, 2018). Consumer activists typically favor regulating AI in general and mandating increasingly transparent XAI in particular<sup>6</sup>. Our paper provides a timely study of the economic implications of regulating XAI on firms and consumers.

We model a duopoly market where firms offer a product based on AI algorithms, e.g., lending based on an AI decision-maker. Firms can set product quality and price while consumers choose to purchase the product from one of the two firms. The policy-maker regulates the AI product by choosing the XAI level ranging from no explanation to full explanations with an objective to maximize the social welfare. In the loan example at full explanations, the AI-based loan rejection decision could be fully explained to the applicant by breaking down how various characteristics of their loan application (credit history, employment, purpose of loan) contribute to their final score. Firms provide this regulated level of XAI through their choice of XAI method. Different XAI methods may explain different characteristics of a rejected loan application while still achieving the same XAI level, e.g., one method may explain employment history while another method may explain how education history contributes to the final score. Firms are thus horizontally differentiated in XAI method and vertically differentiated in quality. The consideration of competitive pricing and

---

<sup>a</sup> By Next Move Strategy Consulting, nextmsc.com

quality allocation alongside XAI adds a much-needed perspective for the ongoing race to mandate full XAI.

As a baseline, we consider a mandatory XAI regulation where the policy-maker chooses the XAI level and firms must offer XAI at this level. We contribute to the philosophical debate on XAI by introducing a new regulatory lever of *optional* XAI where the policy-maker chooses the XAI level, but firms are free to decide whether or not they will offer XAI at this level. One of our first findings confirms the popular belief that mandating XAI usually makes the society as a whole better off. But we also identify conditions where **mandatory XAI has no additional benefit compared to optional XAI**. Intuitively, this occurs when consumers are more sensitive to AI product quality than XAI. As a result, explanations have limited impact on separating demand among firms, resulting in loss of profits. While we do not explicitly model the regulatory cost of monitoring compliance, this finding can lead to significant savings in regulatory oversight while delivering nearly the same outcomes.

A natural follow up question is: What would firms do if left unregulated? At present, firms point to reasons such as long-run implications of stolen IP, copycats, adversarial attacks, or even unavailability of XAI methods (that are widely applicable on any AI model) for not offering full explanations. The computer science literature suggests that these threats and limitations should soon be (if not already are) a non-factor. Therefore, we consciously choose not to model the adversarial threats and set the cost of adopting XAI methods to zero. We do recognize the challenges in organization adoption of AI and XAI strategies, so accordingly we build in that friction into the model. This model of unregulated firm should inform us what the unregulated market will look like as AI expertise matures. We show that an **unregulated market can be as good as regulated** (or even better in settings where firms are limited in price and quality choices) in terms of total welfare, total consumer utility, and average XAI level offered to consumers. The direct impact of XAI is to provide greater utility for consumers and therefore room for all firms to potentially extract at least some of the additional surplus by increasing prices, which explains why optional regulation or unregulated markets work well. Thus, calls for instituting costly mandatory regulations may be redundant.

Closely related to calls for mandatory regulation is the call for full transparency to consumers. It would appear that full explanations are a win-win for firms and consumers. Surprisingly, we actually find that **optimal XAI levels, firm’s choice under no regulation, and policy-maker’s choice under regulation can be less than full explanation**. In fact, full explanation

XAI regulation may make firms and consumers worse off. This result holds for both mandatory and optional XAI regulations and even when the cost of implementing XAI is zero. We also consider an alternative setting where the policy-maker seeks to maximize objectives such as the total number of firms that offer XAI or the average XAI level received by customers. These objectives are more convenient to measure (and publicize) instead of the more comprehensive objective of social welfare. Again, we show that asking for full explanations is not always the best regulatory decision. Further, conventional wisdom would suggest that under optional XAI, since firms have more choices, the policy-maker needs to conservatively set a lower XAI level in order to incentivize firms to opt in to offer XAI. But we show that **requiring a lower XAI level (more relaxed standard) may be the optimal choice under *mandatory* regulation.**

To understand these finding of partial (i.e., less-than-full) explanations, note the indirect impact of XAI in modulating competition between firms. If firms use different XAI methods, their explanations reveal different information sets, catering to different customer segments or groups. At full explanations, the information sets overlap, eliminating any differentiation. This mechanism explains the impetus for partial explanations. This would also suggest that firms may prefer asymmetric XAI levels. It is unclear if consumers or policy-maker would prefer the same. Naturally, a mandatory XAI regulation would not allow any asymmetry. In general, it is difficult for policy-maker to institute asymmetric policy across firms. This is a hidden reason why we examine optional XAI policy, which allows some asymmetry. Without modelling the inter-related economic forces it is difficult to answer – will unregulated firms prefer asymmetry? and should policy-maker allow asymmetry? Policy objective of total welfare doesn't directly say anything about symmetry.

To underline the importance of examining symmetry, let us go back to original motivation for explanations, i.e., trust and adoption of AI. Note that different consumer segments may be *fairly* served (to trust and adopt AI) with different information sets in their explanations. We will discuss how members of one customer group (say protected or historically dis-advantaged) group will be clustered together in terms of the information set that explains their prediction, but distant from other customer group (say historically advantaged). Thus symmetry is closely related to AI fairness. Firms differentiated in XAI methods (thus serving different consumer groups, instead of just the historically dominant group) but offering symmetric XAI attains this fairness objective. In theory a mandatory XAI regulation would ensure this symmetry, but unregulated market may not. Surprisingly, our third finding reveals that even **unregulated firms**

**always choose symmetric XAI levels. In fact, unregulated XAI can lead to greater AI fairness than regulated.** This finding again pushes back on the eagerness to spend on mandatory regulation.

The takeaway message is that policy-makers must consider factors such as firms' XAI methods and the market structure (whether firm competition is dominated by quality or by explanation). They should not race to mandate full explanations, but rather smartly adjust the levers on the XAI level and give options to firms. In terms of managerial implications for firms, we point out the significance of quality as the main AI product attribute when developing XAI strategy. This is true even when firms compete on explanations and consumers value explanations more than quality. There are two intermediate results in our analysis worth highlighting. First, when firms are using the same XAI method, then quality and XAI are substitutes, i.e., the firm that offers higher quality can sacrifice explanation to maintain its differentiation. But, when firms are already differentiated in XAI method, quality and XAI are complements, i.e., the firm with a greater XAI level is also able to invest more in quality due to the market share. Second, equilibrium solutions turn out to be expressed such that firm's marginal cost of quality and consumer utility from explanations are interchangeable. This highlights the underlying trade-off between quality and explanations and the fact that these two must be considered jointly in XAI decisions.

To our knowledge, this is the first attempt towards comprehensively representing AI (quality and accuracy), XAI (level, method, and fairness), and heterogeneous customer preferences for XAI simultaneously. Our model is also generalizable across AI algorithms and XAI methods. It also encompasses multiple regulatory levers (mandatory, optional and no regulation) as well as multiple policy objectives (total welfare, consumer welfare, average explanations and fairness). Tractable closed form solutions for such a comprehensive model is challenging. This is achieved only by cleverly breaking down strategy spaces into regions (e.g., firm competition dominated by quality vs explanation) where welfare and profit maximization is analytically solvable. The thoroughness of the model is not an arbitrary engineering goal, but a conscious decision to model all inter-related economic forces (price, quality, explanations) while assuming away the non-economic forces (e.g., stolen IP and adversarial attacks). This is unique among the literature in this area and hopefully provides a framework for more research questions to be answered.

## 2. Background

XAI methods often use state-of-the-art human-computer interfaces (mostly visual), such as Figure 1. Most of the XAI scholarship focuses on local explanations that justify a specific decision or prediction. These “local” explanations are generated by methods such as LIME, LOCO, and SHAP (Lei et al., 2017; Lundberg and Lee, 2017; Ribeiro et al., 2016). But not all local explanations are equally local: XAI methods can be tuned to reveal different amount of information in the explanations. For example, Facebook provides some local explanation for why a certain ad was shown (Figure 1). Facebook’s ad targeting AI likely uses dozens of customer features, but in this specific example they explain only three features that resulted in the customer receiving a particular ad. In comparison, Google Search page informs the user that a certain ad was shown because of “your current search terms”. This provides only a global description of how the AI works and offers zero local explanation.

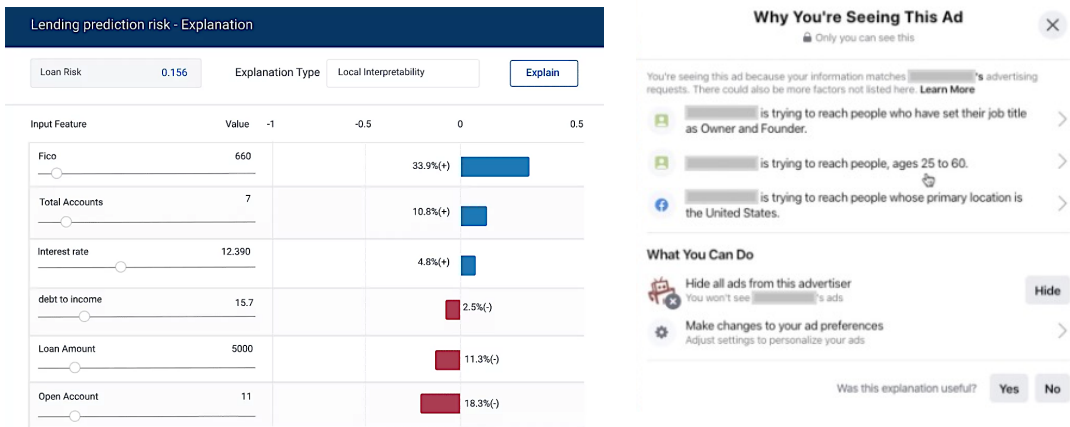


Figure 1. (Left) When applying for a loan, XAI analyzes each feature and finds its positive/negative effect on the outcome (loan risk). Source: AIMultiple<sup>7</sup>. (Right) Facebook provides some explanation about why it shows an ad to you.

### 2.1. The Need for XAI

The term “XAI” was coined in 2004 (van Lent et al., 2004), but the problem of explainability has existed since the work on expert systems in the mid-1970s (Moore and Swartout, 1988). It was only after the proliferation of AI across industries in recent years that social, ethical, and legal pressures began calling for new AI methods that are capable of making decisions explainable and understandable<sup>8</sup> (Adadi and Berrada, 2018). In regard to social justice and fairness, it is well-documented that AI systems might yield biased results

(Caruana et al., 2015; Howard et al., 2017). For example, COMPAS (Lightbourne, 2017, Tan et al., 2018), a black-box recidivism risk assessment software has been criticized for violating the due process rights because it uses gender and race to predict the risk of recidivism. If Judges are to trust the decisions made by such AI systems, some transparency and information is required to ensure that there is an auditable way to prove that decisions made by the algorithm are fair and ethical.

Another area where fair decisions are of highest importance and could benefit from XAI is financial services. In the US, the Fair Credit Reporting Act (FCRA) requires consumer reporting agencies to provide a list of the key factors that negatively influenced the consumer’s score<sup>9</sup> (McEneney and Kaufmann, 2005). Consequently, loan issuers are required by law to make fair decisions in their credit score models and provide the needed “reason code” to borrowers who were denied credit. This has pushed some credit bureaus such as Experian and Equifax to work on AI-based models that generate more explainable and auditor-friendly reason codes<sup>10</sup>.

XAI is also finding its place in some emerging fields such as autonomous vehicles (Bojarski et al., 2016; Haspiel et al., 2018). In fact, a number of recent accidents (some fatal) have fueled societal concerns about the safety of this technology<sup>11</sup> (Stanton et al., 2019; Yurtsever et al., 2020). Thus, the real-time decisions of autonomous vehicles need to be explainable in the sense that they are intelligible by road-users. Current and next-generation intelligent transportation systems must comply with the set of safety standards established using the right of explanation<sup>12</sup> (“Dentons,” 2021).

Finally, recent research in medical diagnosis has focused on making clinical AI-based systems explainable (Ahmad et al., 2018; Caruana et al., 2015; Che et al., 2016; Holzinger et al., 2017a; Katuwal and Chen, 2016). The urgency of XAI in this domain is not just due to trust concerns—it can literally save lives. For example, in the mid-1990s, researchers trained an artificial neural network (ANN) to predict which pneumonia patients should be admitted to hospitals and which should be treated as outpatients. But the AI would not admit pneumonia patients with asthma on the grounds that it thought they have a lower risk of dying, which is both medically incorrect and dangerous (Adadi and Berrada, 2018). It turned out that the reason was that in the training dataset, pneumonia patients with asthma would often get admitted not just to the hospital but sent directly to the ICU and treated intensively. Since most would survive, the AI concluded that pneumonia and asthma together do not increase the risk of death. Only by

interpreting the model using XAI methods can we discover such life-threatening issues and prevent them.

## 2.2. XAI and AI Regulations

The growing body of research on XAI coincides with the emerging issue of the regulatory and policy landscape for AI in jurisdictions across the world (Law Library of Congress (U.S.), 2019). In 2016, the EU introduced a *right to explanation* of algorithmic decisions in General Data Protection Right (GDPR)<sup>a</sup> which gives individuals the right to request an explanation for decisions made by automated decision-making systems that affect the individual significantly, especially legally or financially (Goodman and Flaxman, 2017). GDPR has garnered support both from consumer right groups such as The European Consumer Organization<sup>13</sup> and from firms such as Facebook/Meta who view the law as an opportunity to improve their data management<sup>14,15,16</sup>. There have been calls for GDPR-style laws to be adopted in the US as well<sup>17,18</sup>. FCRA is the closest equivalent of GDPR in the US as of now. Moreover, France, Germany, and the UK operate under a *comply-or-explain* model of corporate governance<sup>19</sup> which obligates private corporations to adhere to a corporate governance code. Should they depart from the code in any way, they must publicly explain their reasons for doing so (Doshi-Velez et al., 2017). In addition, France has amended its administrative code with the Digital Republic Act which creates a right for subjects of algorithmic decision-making by public entities to receive an explanation of the parameters (and their weighting) used in the decision-making process<sup>20</sup>.

However, some critics view AI regulations such as GDPR as unnecessary or even harmful. They argue that regulating AI might stifle many of its social and economic benefits and restrain innovation<sup>21</sup>. Others, such as Intel’s CEO Brian Krzanich, argue that it is too early to regulate AI as it is still in its “infancy”<sup>22</sup>. Critics of GDPR, in particular, highlight the fact that many AI algorithms are not intrinsically explainable. While XAI seeks to alleviate this problem (Miller, 2019; Mittelstadt et al., 2019), some critics have called into question the value of XAI as well. Among them is Peter Norvig, former Google research director who argues that humans are not very good at explaining their decisions either, and that the credibility of the outputs of an AI system could be better evaluated by observing its outputs over time (Adadi and Berrada, 2018). If human decisions

---

<sup>a</sup> See Articles 13 through 15: <https://tinyurl.com/EU-regulation-2016-679>



can depend on intuition or a “gut feeling” that can hardly be put into words, it can be argued that machines should not be expected to meet a higher standard.

### 2.3. Literature

This paper is at the intersection of three growing streams of research. First, a stream of research looks at potentially negative implications of unexplained AI on society (Fu et al., 2022; Malik, 2020) and consumers’ negative perception from unexplained AI decision changes (Bertini and Koenigsberg, 2021). A second stream of research examines concerns with explainable AI arising from tradeoff with AI accuracy (Adadi and Berrada, 2018), privacy of firms’ proprietary secrets (Adadi and Berrada, 2018), potential for adversarial attacks by competitors (Goodfellow et al., 2015), and strategic manipulation by consumers (Wang et al., 2020). A third stream of research examine AI algorithms deployed by competing firms, for example (Assad et al., 2020) study endogenous collusive learning in the AI algorithm. In the next section we lay out our model and contrast with some of this related literature.

## 3. Theoretical Model

We study a duopoly where two firms, indexed  $i = 1, 2$ , are vertically differentiated in product quality and horizontally differentiated in XAI method. Firms market a product based on AI algorithms (e.g., bank loans). The quality of the product  $q_i \geq 0$  can be related to the accuracy of the AI prediction such as the likelihood of default, or it could be unrelated to AI, e.g., timely disbursement of loan or comprehensive paperwork. Without loss of generality (WOLOG), we assume that firm 1 is the high-quality firm and firm 2 the low quality one, that is,  $q_1 \geq q_2$ . A quadratic cost is considered for quality as  $\beta q_i^2$ .

Customers are able to observe product qualities  $q_i$  and prices  $p_i$ . We model customer preferences using a characteristics-based approach, defined directly over the attribute dimensions of the available products. Customers are heterogeneous in terms of their valuation of quality (i.e., willingness to pay). This heterogeneity can be motivated, for example, by difference in income levels and is captured by parameter  $\theta$  which is uniformly distributed over the population. WOLOG,  $\theta$  is normalized to  $[0, 1]$ . Customer  $j$ ’s flow utility derived from quality of firm  $i$  is  $u_{ij}^q := \theta_j q_i$  where  $\theta_j$  is the customer’s sensitivity to quality. Everyone has the same preference for price so that those who purchase from firm  $i$  experience the same price flow (dis)utility  $u_{ij}^p := -p_i$ . Since products are based on AI, the third component of firms’ products is the explanation that they offer as the reasoning

behind the AI decision. The explanation has two characteristics: XAI level  $\xi$  and XAI method  $e$  as follows:

**XAI level  $\xi$**  represents the amount of information contained in the explanation. Consider a loan approval AI that models applicant scores as a function of 8 application features (e.g., credit history, employment, purpose of loan). A high XAI level would explain a loan rejection decision to the applicant by breaking down how, say, 6 out of 8 characteristics of their loan application contribute to their final score. A low XAI level would only explain, e.g., 2 out of 8 characteristics. Instead of number of features explained, one could alternatively interpret XAI level as R-square explained. In an extreme case, zero XAI ( $\xi = 0$ ) would be akin to no explanation.

In our model, customers prefer higher XAI level. But some customer predictions may need more information to be explained or equivalently higher XAI level  $\xi$  than other predictions. Consider again the loan application scoring AI using a decision tree model (Figure 2). The decision tree splits the entire sample iteratively using informative features. Eventually, each branch in the tree terminates because the samples remaining in the terminating node cannot be further split into more informative nodes. Each customer prediction is in exactly one of these terminating nodes. One way to explain the prediction would be to present to the customer all the splits from the starting node to the terminating node and how each split contributed to or deducted from the customers' score. Since some terminating nodes have lower depth, they can be explained by explaining fewer splits. The XAI level  $\xi$  here would refer to the number of splits or features explained. If all splits are explained ( $\xi = 1$ ), then all customer predictions get explained irrespective of their depth. If only some splits are explained (e.g.,  $\xi = 0.5$ ), this may fully explain some predictions at a shallow depth but not others. The same intuition can be derived from a one-dimensional example where some representative points are explained in detail (Figure 3). If a customer's characteristics happen to be at or very close to one of the representative points, they will get a satisfactory explanation for themselves. As XAI level increases, more representative points are explained, and consequently more customers are satisfied by the explanations.

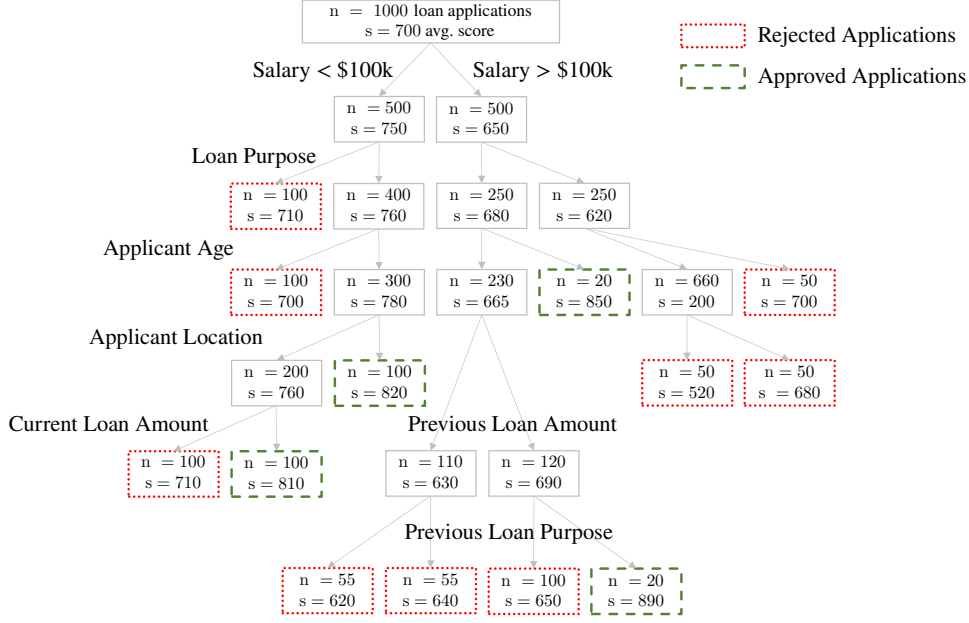


Figure 2. A decision tree that scores loan applications and approves the application if predicted score is  $>800$ . The tree branches terminate when a node has  $\leq 100$  samples. Rejected applications can be explained based on all tree splits. Some rejected applications can be explained using two features, others can be fully explained using 5-6 features.

To model the relationship between customer’s flow utility and explanations, we consider a Hotelling line such that customer AI predictions are uniformly distributed on the  $x$  axis over  $[0, 1]$ . We can think of firm  $i$ ’s XAI as located at a single point on this line at  $e_i$ . At a low XAI level, the firm explains well very few predictions close to  $e_i$ . Customer  $j$ ’s flow (dis)utility due to lack of explanations from firm  $i$  is denoted by  $u_{ij}^e$  and is increasing in the distance between the customer’s preference  $x_j$  and firm  $i$ ’s location  $e_i$  on the Hotelling line, that is,  $u_{ij}^e := -t|x_j - e_i|$  where  $t > 0$  is the “transportation cost” in the Hotelling model and can be thought of as the per unit cost of misfit. An increasing XAI level can be modelled as the firm being located at multiple points on the line such that more predictions are in close vicinity. Alternatively, the same can be modelled as the firm still located at one point  $e_i$  but the transportation cost shrinking—as a function of XAI level  $\xi$ —to  $t(1 - \xi)$ . At full XAI  $\xi = 1$ , the decision tree (Figure 2) explains all splits and the one-dimensional function (Figure 3) will have a large number of explained points (stars), eventually covering all values of  $z$  and explaining all customer predictions.

**XAI method e:** If firms use the exact same XAI method, they are co-located with each other, i.e.,  $e_1 = e_2$ . When both firms offer the same XAI level  $\xi$  (e.g., explaining 10 representative points on the one-dimensional function) but use different XAI methods  $e_1 \neq e_2$  (e.g., explain a different set of 10 points),

explanations offered by each firm better suit a different set of customer predictions. We analytically model two extremes, i.e., same method ( $e_1 = e_2 = 1/2$ , denoted by ‘min’) and diametrically opposite methods ( $e_1 = 0$ ,  $e_2 = 1$ , denoted by ‘max’<sup>a</sup>).

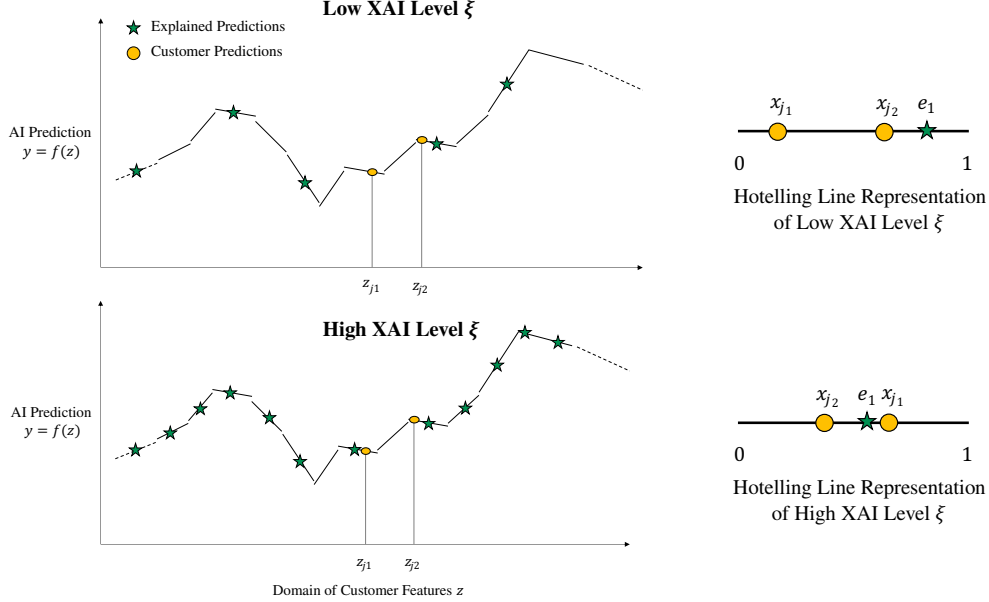


Figure 3. (Left) An AI prediction model  $y = f(z)$ .  $j_1$  and  $j_2$  are two representative customers. Each customer gets a prediction  $y_j$  and an explanation that may not exactly explain their prediction. At low XAI level, available explanation is local for customer  $j_2$  but not for  $j_1$ . At high XAI level, explanation is local for both  $j_1$  and  $j_2$ . (Right) Equivalent Hotelling line representation for low and high XAI levels. All points explained collapse into single firm position  $e_1$  on the Hotelling line while the customer position  $x_j$  is relative to the closest explained point.

Note that the intuitive examples in Figure 2 and Figure 3 are not exhaustive. XAI methods can differ in the set of explained predictions, the set of explained features, the way explanations are presented, or whether the explanations are counterfactual (e.g., your chance of receiving the loan would be more than 50% if your salary were to increase by \$5,000) vs. inferential (e.g., your loan was rejected because your salary was \$4,000 less than average). In practice, consumers can be heterogeneous in their preference for one or the other XAI method based on any of these differences. We consciously do not model specific

<sup>a</sup> An example of different XAI methods is when one firm uses SHAP and the other uses LIME.

XAI methods to ensure that our model covers all current and future XAI methods.

In aggregate, customer  $j$ 's indirect utility from firm  $i$ 's quality, price, and explanation is as follows:

$$u_{ij} = V_j + u_{ij}^q + u_{ij}^p + u_{ij}^e = V_j + \theta_j q_i - p_i - t(1 - \xi_i)|x_j - e_i| \quad 3.1$$

where  $V_j$  is the customer's reservation price or income, i.e., the intrinsic value that he has for the products in this market. We assume a covered market, so  $V_j$  is large enough for all customers that  $u_{ij}$  is always positive in equilibrium and every customer makes exactly one purchase<sup>a</sup>. The customer's utility from explanations  $u_{ij}^e$  depends on the XAI level of firms  $\xi_i$ , the XAI methods  $e_i$ , and customer's preference for explanation  $x_j$ . Next, we elaborate on these as well as related elements of AI accuracy and AI fairness.

**AI Accuracy:** The two attributes of products—quality and explanations—are orthogonal in our model, that is,  $q_i$  and  $\xi_i$  do not necessarily depend on each other. This is true even if product quality  $q_i$  is tied to AI model accuracy. Some researchers have argued that firms may prefer not to use XAI because they will have to sacrifice AI accuracy (Adadi and Berrada, 2018). But more recent research points to the contrary. (Rudin, 2019) argues that such a trade-off between accuracy and interpretability is a “blind belief” and a “myth” that has no real data to support it. Rudin shows that in problems with structured data and meaningful features, there is often no significant difference in accuracy between more complex classifiers (deep neural networks, boosted decision trees, random forests) and much simpler classifiers (logistic regression, decision lists) after pre-processing. Therefore, our model consciously does not incorporate any trade-off between XAI and AI accuracy (as one of the measures of quality).

**AI Fairness:** The representation described above implicitly captures this important notion of AI fairness. Note that motivation of XAI is to encourage adoption and trust. If XAI in equilibrium turns out to be less than full XAI i.e.,  $\xi < 1$ , then different customer predictions on the hoteling line  $x_j$  receive unequal explanations. Consider two customer groups – protected (historically marginalized) and unprotected. The predictions for the two groups may be best explained by different features e.g., employment or education for one group but

---

<sup>a</sup> This, of course, makes use of the common assumption in the discrete-choice framework literature that income effects from price changes are negligible (Cunha et al., 2020), that is, income and prices are additive separable. Therefore,  $V_j$  can be omitted from Eq. 3.1 since it does not vary across products.

credit history and purpose of loan for another group. On the hoteling line, members of protected group will be closely clustered together (say on right extreme) but distant from the unprotected group (say in the middle). If the two firms use the same XAI method (situated in the middle) and offer less than full explanations, the protected group systematically receive less explanation. In theory, such an unequal outcome is alleviated by mandatory XAI regulation (full or close to full) and differentiated (“max”) XAI methods. The latter helps because the firms when revealing different information in their explanations, from the virtue of different XAI method, indirectly cater to different customer groups. There is a natural concern – would unregulated firms choose asymmetric XAI levels (potentially to weaken competition), thus result in unfair XAI?

Moreover, some research looks at firm’s tradeoff between XAI and protecting the privacy of its AI model and trade secrets (Adadi and Berrada, 2018). It is shown that it may be possible to “steal” underlying AI models if AI input and outputs are provided (Barredo Arrieta et al., 2020; Orekondy et al., 2019). AI algorithms can also be subject to adversarial attacks that aim to confuse the model to lead it to a desired output by an adversary (Goodfellow et al., 2015). One could extrapolate from this and argue that firms may prefer not to use XAI because it may foster such issues. We are not aware of research showing that XAI has resulted in loss of a firms’ intangible assets, led to copycats, or adversarial attacks. Thus, we exclude these factors from our model when considering firms’ choice of XAI. In addition, recent research argues that in some situations, firms may not use provide transparency to avoid strategic manipulation by consumers who learn to game the AI decision-maker over repeat interactions (Wang et al., 2020). We consciously exclude consumer learning over time as a factor.

### 3.1. Market Structures

According to above,  $(e_1, q_1)$  and  $(e_2, q_2)$  represent the positions of firm 1 and firm 2 on the  $q$ - $x$  plane, respectively. At the same time, customers’ preferences for quality and explanation are uniformly distributed over a unit square on the  $\theta$ - $x$  plane such that customer  $j$  is represented by  $(x, \theta) \in [0, 1]^2$ . Firm  $i$ ’s demand (market share)  $d_i$  is the area of this unit square that firm  $i$  captures. Obviously,  $0 \leq d_i \leq 1$  and  $d_1 + d_2 = 1$ .

For any customer of type  $(x_j, \theta_j)$ , the utility from the product offered by firm 1 and firm 2 is  $u_{1j}$  and  $u_{2j}$ , respectively, where<sup>a</sup>:

---

<sup>a</sup>  $V_j$  is omitted due to the previous footnote.

$$\begin{aligned}
u_{1j} &= \theta_j q_1 - t(1 - \xi) |x_j - e_1| - p_1; \\
u_{2j} &= \theta_j q_2 - t(1 - \xi) |x_j - e_2| - p_2.
\end{aligned}
\tag{3.2}$$

Firm 1's demand  $d_1$  are customers for whom  $u_{1j} > u_{2j}$ . Similarly, firm 2's demand  $d_2$  are those for whom  $u_{2j} > u_{1j}$ . The set of customers who are indifferent between making a purchase from firms 1 and 2 is called the *marginal customers* and is represented by a line dividing the unit square on the  $x$ - $\theta$  plane. This is called the *indifference line* and it is found by setting  $u_{1j} = u_{2j}$ . Depending on the location of firms on the  $x$  axis and the value of  $\xi_1$  and  $\xi_2$ , the slope and intercept of the indifference line will be different, which in turn has an impact on equilibrium results. Following (Neven and Thisse, 1989; Vandenbosch and Weinberg, 1995; Wattal et al., 2009), we analyze the possible markets formed by the indifference line separately.

**Lemma 1. Existence of Two Market Structures**

*Depending on the slope and the intercept of the indifference line, only two market structures can have Nash equilibrium: the “explanation-dominated” market and the “quality-dominated” market, denoted by  $E$  and  $Q$ , respectively.*

✂ Proof in Appendix A1.

In an explanation-dominated market (or **market E** for short), what separates customers is their taste in explanations rather than quality (Figure 4, left). Since  $q_1 \geq q_2$ , customers on the left side of the indifference line are firm 1's demand ( $d_1$ ) and customers on the right are firm 2's ( $d_2$ ). Here, for any  $\theta_j \in [0, 1]$  there exists a customer  $(x_j, \theta_j)$  who is indifferent about the firms' products. By contrast, in a quality-dominated market (or **market Q**), customers are separated based on their taste in quality (Figure 4, right). Since  $q_1 \geq q_2$ , customers above the indifference line are firm 1's demand ( $d_1$ ) and below it firm 2's ( $d_2$ ). In this market, for any  $x_j \in [0, 1]$  there exists a customer  $(x_j, \theta_j)$  who is indifferent about the firms' products.

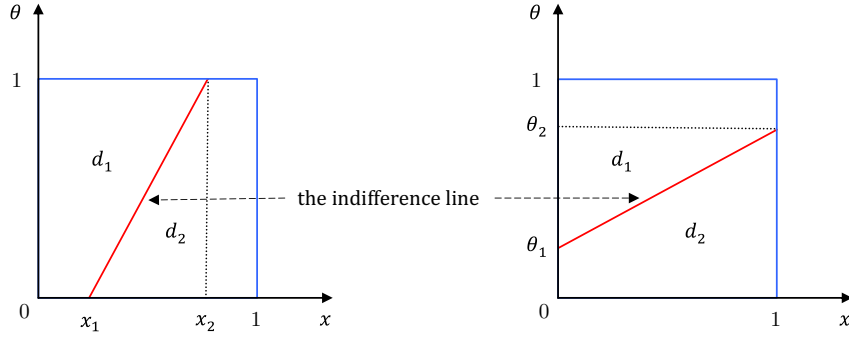


Figure 4. (Left) Explanation-dominated market. (Right) Quality-dominated market

### 3.2. The Policy-Maker's Problem

The policy-maker's goal is to set an XAI level  $\xi$  for both firms such that the *total welfare*  $W_t$  is maximized.  $W_t$  is composed of firm profits and consumer utility, which is calculated by integrating Eq. 3.1 over  $d_1$  and  $d_2$ .

We consider two policy setups: (1) **Optional** XAI, (2) **Mandatory** XAI. In mandatory XAI, the policy-maker sets the XAI level  $\xi$  and enforces this policy. To do that, the policy-maker calculates firm profits and consumer utility—knowing that both firms offer XAI—and finds the optimal  $\xi$ . In optional XAI, the policy-maker sets the XAI level  $\xi$  as a guideline, but firms are free to choose whether they will offer XAI at this level or they will offer no XAI at all. The policy-maker must predict firms' equilibrium strategies and set optimal  $\xi$  accordingly. Through the choice of  $\xi$ , the policy-maker can shape the market structure, and hence firms' equilibrium strategies. We will also discuss scenarios where for some values of  $\xi$ , both market structures E and Q are possible at the same time. Firms will choose the market that yields higher profits, but in case they have conflicting interests, there will be no equilibrium. The policy-maker must avoid this situation by choosing  $\xi$  properly.

### 3.3. The Firms' Problems

Firm  $i$ 's problem is to maximize its profit function  $\pi_i$  which is defined as:

$$\pi_i = p_i d_i - \beta q_i^2. \quad 3.3$$

Firms face two situations:

- Regulated XAI level: The policy-maker sets  $\xi$  for both firms.
- Unregulated XAI level: Firms freely choose their own  $\xi_i$ .

Under regulated XAI level, the policy-maker sets  $\xi_1 = \xi_2 = \xi$ , so  $\xi$  is endogenous for the policy-maker and exogenous for firms. In this situation, each



firm  $i$  has to maximize Eq. 3.3 by choosing the optimal  $q_i^*$  and  $p_i^*$ . The third decision that firms must make has to do with offering XAI. As mentioned before, under optional XAI, each firm must choose whether or not it offers XAI at level  $\xi$  set by the policy-maker. But under mandatory XAI, both firms are forced to offer XAI at level  $\xi$ .

Under unregulated XAI level, there is no policy-maker involved and each firm  $i$  maximizes Eq. 3.3 by finding the optimal  $\xi_i^*$ ,  $q_i^*$ , and  $p_i^*$ . This setting is not a primary focus of our paper and is only briefly discussed in Section 6 because we foresee the policy-maker taking an active role, even if for political or PR reasons.

### 3.4. The Games

Firm  $i$ 's choice of XAI method which determines location  $e_i$  on the Hotelling line is a long-term strategy. It often needs multiple levels of approvals within the organization and with the regulators. The way XAI is presented also shapes how the customers interact with the firm's product. Thus, changing  $e_i$  requires significant investment in product repositioning. As a result, we treat firms' locations on the Hotelling line as exogenous decisions made before the start of the game. But while switching XAI method is too costly, implementing XAI is assumed costless in our model. The reason is that explanations are now readily available if a firm adopts methods such as SHAP or LIME. These model-agnostic methods provide explanations out of the box at no extra cost.

We discuss the necessary setup for game solution for the optional XAI regulation where the policy-maker sets XAI level  $\xi$ . The game solutions for mandatory XAI will be simplified version of the same. There are two market structures E and Q under two firm differentiations 'max' and 'min'. So, there are 4 games to analyze. Table 1 illustrates 3 equilibrium scenarios mentioned in Section 3.2 for each of these 4 games.

Table 1. Labels for 4 games and 3 equilibrium scenarios that could happen in each game.

		Firms' Horizontal Differentiation			
		max		min	
		Market Structure		Market Structure	
		E	Q	E	Q
How Many Firms Offer XAI?	0	maxE0	maxQ0	minE0	minQ0
	1	maxE1	maxQ1	minE1	minQ1
	2	maxE2	maxQ2	minE2	minQ2
		Game 1	Game 2	Game 3	Game 4

Each game has three sequential stages:

- **Stage 0:** Firms simultaneously decide whether they will offer XAI  $\xi_i \in \{\xi, 0\}$ ;
- **Stage 1:** Firms simultaneously choose their quality levels  $q_i$ ;
- **Stage 2:** Firms simultaneously choose their prices  $p_i$ .

The intuition behind the three-stage game structure is the fact that prices are more flexible than quality in the short-term (Cunha et al., 2020). Likewise, product quality (e.g., through the AI model's accuracy) is more flexible than XAI strategy. Thus, decisions in early stages can be viewed as the firm's long-term strategy while subsequent stages involve progressively short-term decisions.

Stage-0 equilibrium decisions could be one of the following:

- No firm offers XAI:  $(\xi_1^*, \xi_2^*) = (0, 0)$
- One firm offers XAI:  $(\xi_1^*, \xi_2^*) = (\xi, 0)$  or  $(0, \xi)$
- Both firms offer XAI:  $(\xi_1^*, \xi_2^*) = (\xi, \xi)$

For example, with  $(\xi_1^*, \xi_2^*) = (\xi, 0)$  firm 1 offers explanations at equilibrium while firm 2 does not. At each game, we use backward induction to solve for equilibrium qualities and prices in stages 1 and 2, respectively.

## 4. Regulated or Exogenous XAI

### 4.1. Optional XAI

First, we present a summary of stage-0 equilibrium decisions in Table 2. Then we discuss the equilibria in markets Q and E. Having found the equilibria, we will then address the problem of policy-making for optional XAI in Section 4.1.3.

Table 2. Stage-0 equilibrium decisions  $(\xi_1^*, \xi_2^*)^a$ 

		Market Structure	
		E	Quality-Dominated (Q)
Horizontal Differentiation	max	$\begin{cases} (\xi, 0) & \xi > \xi_* \\ (0, 0) & \text{else} \end{cases}$	$(\xi, \xi)$
	min	$(\xi, 0)$	$(\xi, \xi)$

#### 4.1.1. Equilibria in Quality-Dominated Markets (Q)

As we see in Table 2, the quality-dominated market is the only market structure in which both firms offer explanations at level  $\xi$  set by the policy-maker, and this is irrespective of firms' horizontal differentiation (max and min). Moreover, we find that  $(\xi_1^*, \xi_2^*) = (\xi, \xi)$  *weakly* dominates  $(\xi_1^*, \xi_2^*) = (0, 0)$  (See Table 3). Intuitively, this is due to symmetric  $\xi$  and the fact that firms capture the entire Hotelling line in market Q (See Figure 5). Therefore, the misfit cost  $t$  has no effect on the equilibrium<sup>b</sup>.

Table 3. 4 identical scenarios in quality-dominated markets

maxQ0	maxQ2	minQ0	minQ2
$p_1^* = \frac{4}{27\beta}, \quad p_2^* = \frac{2}{27\beta}, \quad q_1^* = \frac{2}{9\beta}, \quad q_2^* = 0, \quad \pi_1 = \frac{4}{81\beta}, \quad \pi_2 = \frac{2}{81\beta}$			

---

<sup>a</sup> For the condition of each equilibrium, see Appendix A4.

<sup>b</sup> One might consider a *fixed cost*  $1\{\xi > 0\} \cdot c_{\text{XAI}}$  for XAI, say, due to implementing XAI for the first time. However, this would not change equilibrium results because Table 3 tells us that equilibrium price, quality, and profit are independent of  $\xi$ . In other words,  $(\xi_1^*, \xi_2^*) = (\xi, \xi)$  still weakly dominates  $(\xi_1^*, \xi_2^*) = (0, 0)$  in market Q.

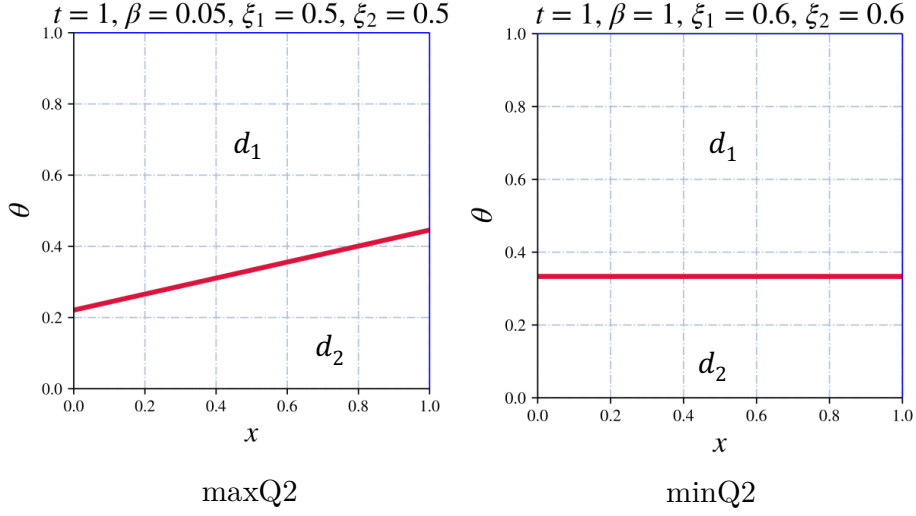


Figure 5. Equilibria, indifference lines, and demands in the Q market for some values of  $\beta$  and  $t$

#### 4.1.2. Explanation-Dominated Market (E)

When firms use different XAI methods, we find that no firm has the incentive to offer explanations at low levels ( $\xi < \xi_*$  where  $\xi_* = 1 - 1/36\beta t$ ). That said, if the policy-maker sets  $\xi > \xi_*$ , then firm 1 (the high-quality firm) provides more explanations while also offering a *higher* level of quality compared to a baseline in which no firm offers explanations. That is, XAI and quality are complements in this case. Firm 2 (the low-quality firm) responds by lowering its quality. Therefore, on one hand we have a quality gap which also grows in  $\xi$  and contributes to greater profits. And on the other hand, explanations reduce the horizontal differentiation between the two firms on the Hotelling line, which should lower their profits.

The net effect of these two forces is as follows: As long as firm 1 offers XAI at high levels ( $\xi > \xi_*$ ), its profit increases with  $\xi$  while firm 2's profit declines. The intuition is that as  $\xi$  increases, firm 1 increases and firm 2 decreases its quality level, such that firm 1 always charges a higher price than firm 2. Therefore, firm 1 captures more customers while also increasing its price, while firm 2 loses its market share while lowering its price (See Figure 6, middle). We show that with firm 1 offering XAI, firm 2 cannot do any better by adopting XAI, ruling out the  $(\xi_1^*, \xi_2^*) = (\xi, \xi)$  strategy.

The situation is different when firms use the same XAI method. As seen in Table 2, firm 1 always finds it more profitable to offer explanations while firm 2 does not. Interestingly, firm 1 (the high-quality firm) *lowers* its quality as  $\xi$  increases while firm 2 (the low-quality firm) does the opposite. In other words, XAI and quality are *substitutes* in this case. We also find that the profits of both

firms increase as more explanations are asked by the policy-maker ( $\partial\pi_i/\partial\xi > 0, i = 1, 2$ ). The reason is as follows:

On one hand, profits increase with quality for both firms ( $\partial\pi_i/\partial q_i > 0, i = 1, 2$ ), leading them to choose almost identical qualities if the cost of quality approaches zero<sup>a</sup>, hence a reduction in the quality gap and profits. On the other hand, offering XAI introduces a measure of horizontal differentiation even though firms are located at the same point on the Hotelling line. This is because the increase in customer utility due to explanations is higher for customers who are farther away from the center than for those closer (See Figure 6, right). This is shown below:

$$\begin{aligned} u_1^e &= -tx(1 - \xi); \\ u_2^e &= -tx; \\ \Rightarrow \Delta u^e &:= u_1^e - u_2^e = t\xi x. \end{aligned} \tag{4.1}$$

As  $x$  gets larger, so does  $\Delta u^e$ , which enables the firm that offers explanations (firm 1) to capture customers located away from the center, while the other firm (firm 2) captures customers located close to the center. Thus, although firms are undifferentiated, XAI by one firm helps both firms earn higher profits.

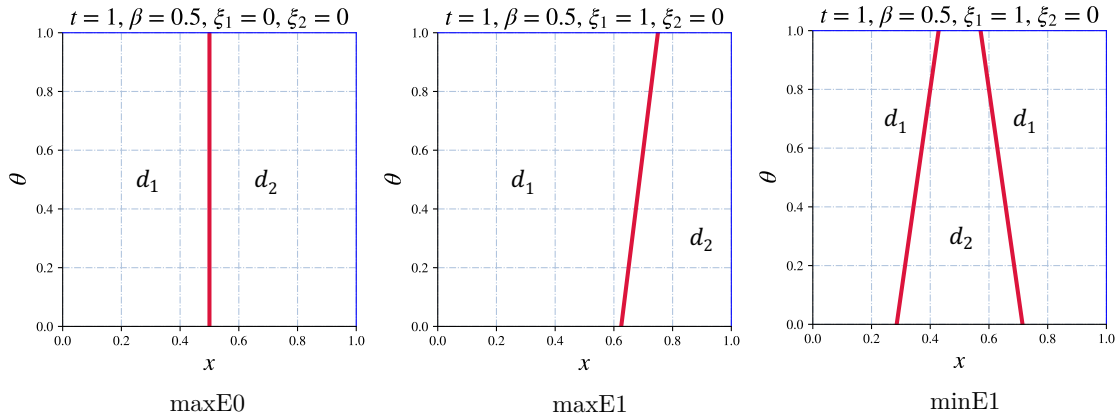


Figure 6. Equilibria, indifference lines, and demands in the E market<sup>b</sup> for  $\beta = 0.5$  and  $t = 1$

#### 4.1.3. Policy-Making for Optional XAI

The policy-maker aims to maximize the total welfare by choosing the optimal  $\xi = \xi_{\text{policy}}^{\text{opt.}}$  for each value of  $\beta t$ , where  $\xi_{\text{policy}}^{\text{opt.}}$  is the policy under optional XAI. Since

<sup>a</sup> Which is not possible because of the equilibrium condition ( $\beta t \xi > 2/9$ ).

<sup>b</sup> Notice that due to the symmetry in minE1, the indifference line in  $x \in [0, 0.5]$  will be reflected in  $x \in [0.5, 1]$  as well, hence the two lines shown in Figure 6.

firms' locations on the Hotelling line is exogenous, the policy-maker must devise an  $\xi$  for 'max' and 'min' separately. The first step is to calculate  $W_t$ . If only one market (E or Q) is possible, then the policy-maker maximizes  $W_t$  in that market w.r.t.  $\xi$ . But if two markets are possible at the same time, the strategic policy-maker calculates  $W_t$  for each market and then chooses  $\xi$  such that firms only choose his desired market and maximum total welfare is obtained.

Table 7 in Appendix A4 lists all equilibrium results and their conditions. According to the condition of each equilibrium, the  $\beta t$ - $\theta$  plane is divided into several regions as shown in Figure 7. While we have analytically derived  $W_t$  for all equilibria<sup>a</sup>, in the cases of maxE1, minQ2, and minE1 it is overly complex, making it analytically intractable to compare  $W_t$  in different regions and calculate the optimal  $\xi$ . To tackle this problem, we run Monte-Carlo simulations on 40,000 points in the unit square in Figure 7. In doing so, we also find the markets that firms choose to form, and areas where firms have conflicting incentives (shown in gray  $\otimes$ ). In these regions, one firm prefers market E and the other prefers market Q, leading to no Nash equilibrium in the market.

In Figure 7, maximum total welfare is achieved at the points on the (thick) blue lines. Contrary to expectation, these graphs reveal that **policy-makers should not always require full explanations even under optional XAI**. This result—which is also irrespective of firms' differentiation—is because for some values of  $\beta t$ , high XAI level leaves the firms indecisive about the market they seek to form. This can be verified analytically for the 'min' case:

**Lemma 2.**

*When firms use the same XAI method and both markets E and Q are possible (Figure 7, right, gray area  $\otimes$ ), the high-quality firm always earns more in market E than it does in market Q, but the opposite is true for the low-quality firm. Choosing  $\xi_{\text{policy}}^{\text{opt.}}$  in this region, then, leads to no Nash equilibrium.*

✖ Proof in Appendix A6.

---

<sup>a</sup> See Appendix A7.

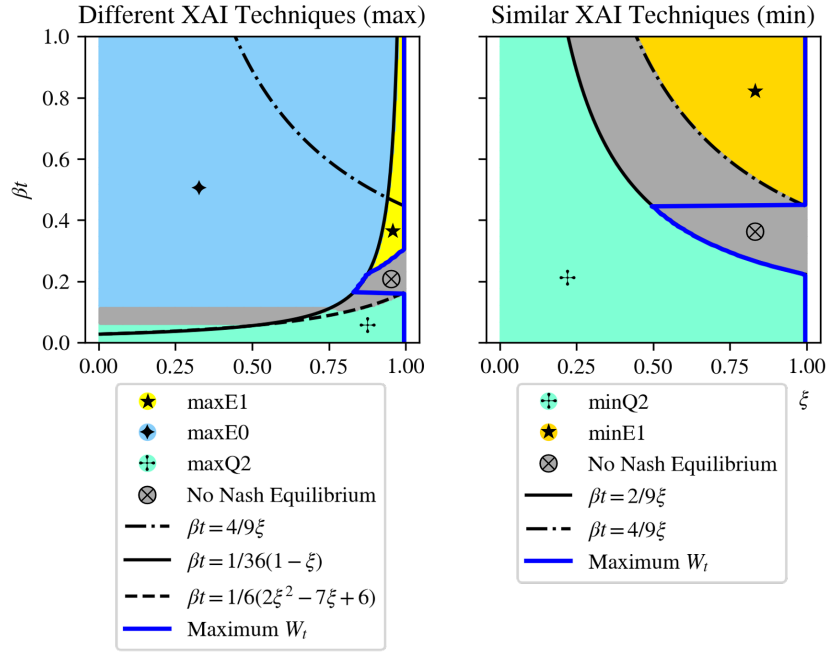


Figure 7. Different regions on the  $\beta t$ - $\xi$  plane and the equilibria that they support under optional XAI

## 4.2. Mandatory XAI

If it is mandatory for firms to offer XAI at level  $\xi$  set by the policy-maker, Table 1 reduces to only 4 scenarios shown in Table 4.

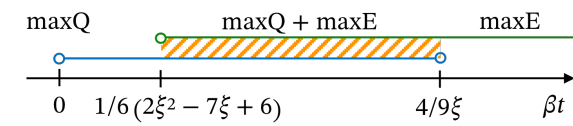
Table 4. The scenarios in which both firms offer XAI

	Scenario			
	maxE2	maxQ2	minE2	minQ2
Is it equilibrium in optional XAI?	No	Yes (in maxQ)	No	Yes (in minQ)

As we saw in Section 4.1.1, maxQ2 and minQ2 in Table 4 are equilibrium in quality-dominated markets. This implies that even without policy enforcement, both firms offer XAI in these scenarios, regardless of the  $\xi$  level. Moreover, scenario minE2 is discarded because the indifference line is horizontal and firms engage in a price war while racing to the bottom of quality, rendering this scenario impossible to sustain. So, there is only scenario maxE2 that is possible but is not equilibrium by default.

The condition for each scenario, total welfare  $W_t$ , and its derivative w.r.t.  $\xi$  are presented in Table 5. As we can see, unlike the case of optional XAI, here comparing  $W_t$  expressions is analytically tractable and there is no need for Monte-Carlo simulations.

Table 5. The 3 possible scenarios under mandatory XAI

	maxE2	maxQ2	minQ2
$W_t$	$\frac{t(\xi - 1)}{4} + \frac{1}{36\beta}$	$\frac{t(\xi - 1)}{2} + \frac{1}{27\beta}$	$\frac{t(\xi - 1)}{4} + \frac{1}{27\beta}$
$\partial W_t / \partial \xi$	$t/4 > 0$	$t/2 > 0$	$t/4 > 0$
Condition	$\beta t > \frac{1}{6(2\xi^2 - 7\xi + 6)}$	$\beta t < \frac{4}{9\xi}$	$\beta t < \frac{4}{9\xi}$
			

In Table 5, the greatest total welfare among all scenarios belongs to maxQ2, followed by minQ2 and then maxE2, and this holds for all  $\xi \in [0, 1]$ . As mentioned above, when firms use the same XAI method ('min'), only scenario minQ2 is possible. As such, the policy-maker sets

$$\xi_{\text{policy}}^{\text{mand.}} = \min\{4/9\beta t, 1\} \quad 4.2$$

to maximize  $W_t$  in this scenario where  $\xi_{\text{policy}}^{\text{mand.}}$  is the XAI level chosen by the policy-maker under mandatory XAI.

On the other hand, two scenarios are possible when firms use different XAI methods ('max'). Looking at the conditions of maxE2 and maxQ2, we find that both scenarios are possible for a range of  $\beta t$ . In this interval, market Q yields a higher total welfare than market E and is thus preferable by the policy-maker<sup>a</sup>. But firms select the market that maximizes their profit, which may not necessarily maximize  $W_t$ . Therefore, the policy-maker must choose  $\xi_{\text{policy}}^{\text{mand.}}$  such that only one market is possible and  $W_t$  is maximized.

<sup>a</sup> See the expressions of  $W_t$  in Table 5.



Based on the conditions of each market, the  $\beta t$ - $\xi$  plane is split into several regions as illustrated in Figure 8 (left)<sup>a</sup>. Based on Eq. 4.2 and the policy that guarantees maximum  $W_t$  in Figure 8 (right), we maintain that **mandating firms to offer full explanations may actually make everyone worse off**, including the very consumers that the policy aims to support.

This result—which is analogous to the case of optional XAI—holds regardless of firms’ horizontal differentiation (‘max’ or ‘min’). In the case of ‘min’, it is because there will be no equilibrium if  $\xi_{\text{policy}}^{\text{mand.}} > 4/9\beta t$ . And in the case of ‘max’, it is because with  $\xi_{\text{policy}}^{\text{mand.}} = 1$ , firms may end up in the red region ① in Figure 8 (left) in which  $\pi_1, \pi_2 < 0$  and firms leave the market, resulting in  $W_t = 0$ .

A question worth asking at this point is: Why is  $\xi_{\text{policy}}^{\text{mand.}} < 1$  for large  $\beta t$ ? The answer requires an explanation of the red region ① in Figure 8, and it is as follows<sup>b</sup>: We know that  $(\xi_1^*, \xi_2^*) = (\xi, \xi)$  is not the equilibrium strategy of maxE by default<sup>c</sup>. Mandating XAI, then, makes both firms worse off compared to their equilibrium strategy. The intuition is that under the mandatory regulation, as  $\xi$  increases, the horizontal differentiation between the firms on the Hotelling line decreases because each firm’s locational differentiation to the other is  $t(1 - \xi)|e_1 - e_2|$ . This in turn leads to price competition and lower profits. Worse yet, there is no quality gap between the firms to increase the profits ( $\Delta q = 0$ ). The result is a decline in the profits as shown in Figure 9.

In Figure 8 (right),  $\xi_{\text{policy}}^{\text{mand.}}$  converges to 1 as  $\beta t \rightarrow \infty$ , suggesting that when the cost of quality and misfit cost are high enough, the policy-maker can require more explanations, but she can ask for full explanations only when  $\beta t < 4/9 = 0.\bar{4}$ . Overall, our results in this Section and before point to the significance of  $\beta$  (hence the marginal cost of quality) and  $t$  (hence the misfit cost) in devising XAI policies, both in optional and mandatory XAI.

---

<sup>a</sup> See Appendix A8 for more details and mathematical expressions of regions and curves.

<sup>b</sup> See also Appendix A8.

<sup>c</sup> In Section 4.1.2, we found the equilibrium strategies to be  $(\xi_1^*, \xi_2^*) = (0, 0)$  for  $\xi < \xi_*$  and  $(\xi_1^*, \xi_2^*) = (\xi, 0)$  for  $\xi > \xi_*$ .

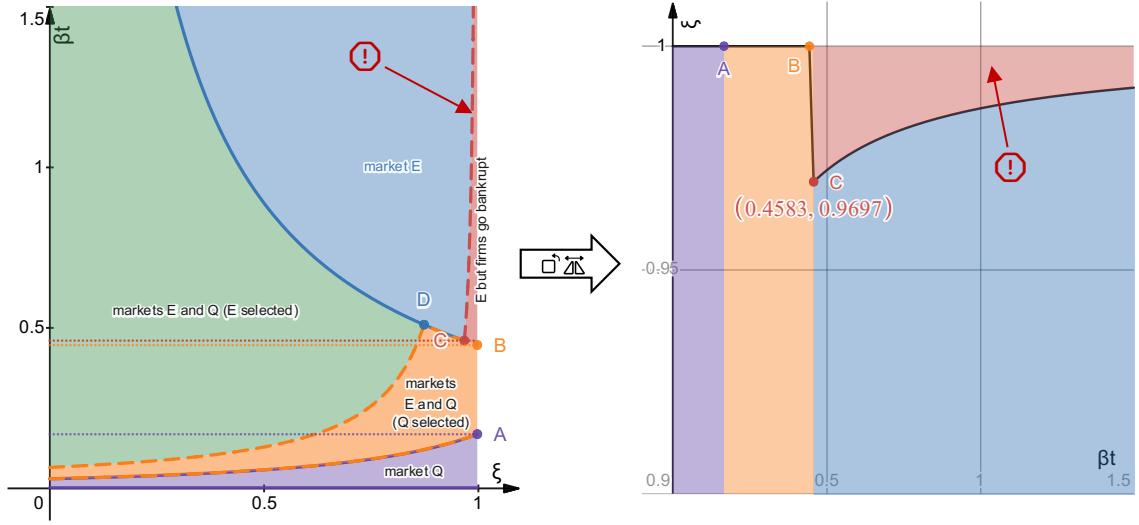


Figure 8. (Left) Different regions on the  $\beta t$ - $\xi$  plane and the markets that they support under mandatory (Right) The policy  $\xi_{\text{policy}}^{\text{mand.}}$  that guarantees maximum total welfare

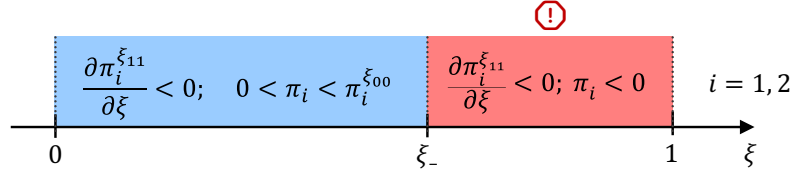


Figure 9. Both firms will be worse off if they use different XAI methods and are mandated to offer XAI in market E

## 5. Regulatory Insights

### 5.1. Should XAI be Mandatory?

In light of the results of Section 4, a natural question is which one—optional or mandatory XAI—leads to higher total welfare? The answer may not be obvious at first glance. The total welfare does not simply depend on the explanation level; it also depends on the number of firms who offer XAI as well as price and quality levels. Take the case of ‘max’ under high  $\beta t$  for instance. On one hand, under optional XAI, the policy-maker asks for full explanations ( $\xi_{\text{policy}}^{\text{opt.}} = 1$ ) but *only one* firm will choose to offer XAI (see Figure 7). On the other hand, if XAI is mandatory, the policy-maker asks for less than full explanations ( $\xi_{\text{policy}}^{\text{mand.}} < 1$ ) but *both* firms must comply (see Figure 8). Thus, one setting has higher explanation level while the other has more number of firms who offer XAI. The net effect is illustrated in Figure 10.

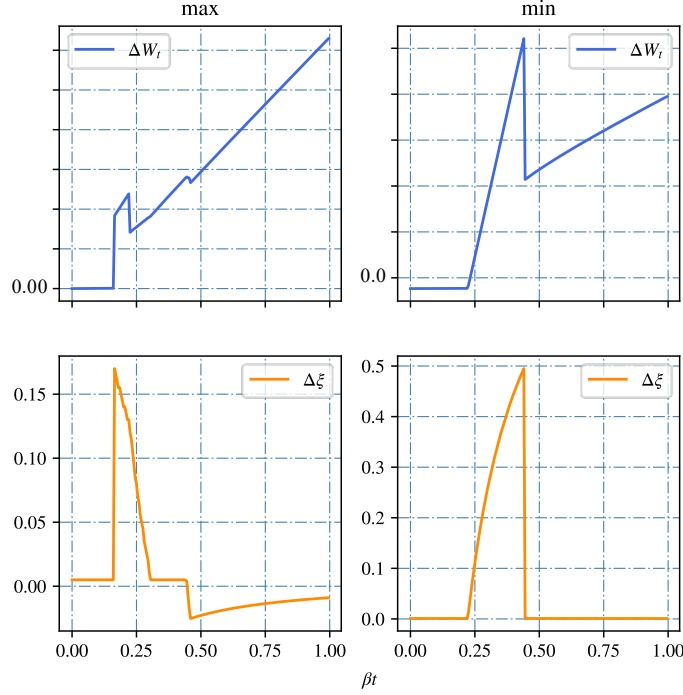


Figure 10. The gaps between total welfare and explanation levels of mandatory and optional XAI.  $\Delta W_t := W_t^{\text{mand.}} - W_t^{\text{opt.}}$ , and  $\Delta \xi := \xi_{\text{policy}}^{\text{mand.}} - \xi_{\text{policy}}^{\text{opt.}}$ .

First, our findings support the intuition that in general, mandating XAI policies makes the society as a whole better off ( $\Delta W_t \geq 0$ ). That said, **we find that there is no additional benefit from mandating XAI ( $\Delta W_t = 0$ ) when  $\beta t$  is small.** The intuition is that small cost of misfit  $t$  suggests that the market does not care so much about explanations, and demands are separated based on customers' taste in quality rather than explanations (quality-dominated market). As such, explanation levels have no effect on firms' profits, even when  $\xi$  is 1. Contrary to conventional wisdom, we also find out that **full explanations are not always necessary to achieve maximum total welfare.** In fact, mandatory policies may actually require *less* explanation levels than optional policies (Figure 10, bottom left), because while optional XAI level can be 1 for  $\beta t > 4/9 = 0.\bar{4}$ , mandatory XAI level must remain less than 1 to allow for market equilibrium<sup>a</sup>.

## 5.2. What if the Policy-Maker's Objective is Different?

In Section 3.2 we mentioned that the policy-maker's goal is to maximize the total welfare of the society. But in practice the policy-maker may find it

<sup>a</sup> As mentioned at the end of Section 4.2.

challenging to measure the total welfare. For political or PR reasons, the policy-maker may seek to maximize an objective that is easier to measure and publicize, such as the total *number of firms* that opt-in to offer XAI or the *average level of explanations* received by customers.

Let us consider the objective that aims to maximize the total number of firms who offer XAI. According to Figure 11 (left), when firms have the same XAI method, the number of firms that offer XAI is *not* monotonically increasing in  $\xi$ . Therefore, **asking for full explanations is not always the best regulatory decision**<sup>a</sup>. This also holds if the objective is the average XAI level  $\xi_{\text{avg.}}$  received by customers<sup>b</sup>, because  $\xi_{\text{avg.}}$  is *not* always monotonically increasing in  $\xi$  (e.g., in Figure 11, left). Figure 11 also reveals that firms' XAI methods influence the trend: If firms use similar XAI methods, the number of firms that offer XAI is, at best, constant, and at worst, decreasing in  $\xi$ . But if firms use different XAI methods, the number of firms that offer XAI is, at worst, constant, and at best, increasing in  $\xi$ .

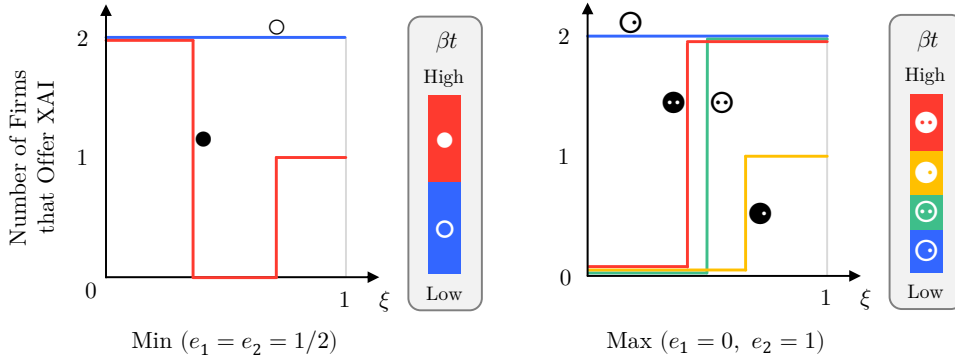


Figure 11. The number of firms that offer XAI as a function of  $\xi$ . Each line corresponds to a different  $\beta t$  level.

## 6. Unregulated XAI

Until now, each firm's XAI decision was either a binary option of  $\{0, \xi\}$  (in optional XAI) or fixed to the policy-maker's choice of XAI level (in mandatory XAI). A natural question is: What would firms do if left unregulated? In our model, this would correspond to firms' choices of XAI levels  $(\xi_1, \xi_2)$ . To answer this, we extend the method in (Wattal et al., 2009) by deriving equilibrium

<sup>a</sup> See Appendix A9.

<sup>b</sup>  $\xi_{\text{avg.}} \stackrel{\text{def}}{=} d_1 \xi_1 + d_2 \xi_2$

results when  $\xi$  is endogenous for firms and not necessarily symmetric<sup>a</sup>. One interesting finding is that when firms adopt different XAI methods in market E, an *asymmetric* equilibrium is not sustainable. In particular, firms always mirror each other's XAI levels<sup>b</sup> such that  $\xi_1 = \xi_2 = 1 - 1/36\beta t$ . This implies that  $\xi_1, \xi_2$  will never be 1<sup>c</sup>. Markedly, this result indicates that **firms choose to hide some information in their XAI outputs even though XAI does not cost anything**. As the misfit cost  $t$  increases and the market cares more about explanations,  $\xi_1, \xi_2$  approach 1. Unfortunately, analytical derivations for market Q are expectedly more complex than the regulated case presented in Section 4; we are only able to attain closed form solutions for the E-market.

To tackle this problem, we use a simplified model as follows: Firms have a binary choice of quality  $\{0, q\}$  and price  $\{0, p\}$  instead of the continuous choice in the full model. Moreover, firms continue to have the full choice of XAI levels, i.e.,  $\xi_1, \xi_2 \in [0, 1]$ . In this simplified model, once again we find that market E is only viable when firms are differentiated in XAI method. Besides, firms' choice of XAI levels is always symmetric,  $\xi_1^* = \xi_2^*$ , and can be less than 1. This is all consistent with the full model. But now using the simple model we can examine the Q-market and once we have the full equilibrium characterization, we can compare the unregulated and regulated settings<sup>d</sup>. Table 6 summarizes the equilibrium XAI levels of our simplified model:

---

<sup>a</sup> We have provided more detailed results in Appendix A10.

<sup>b</sup> But while mirroring each other's XAI *levels* leads to a symmetric equilibrium in this market, mirroring each other's XAI *methods* results in no equilibrium at all. See Appendix A10.

<sup>c</sup> It turns out that  $\xi = 1 - 1/36\beta t$  is the same as the line  $\beta t = 1/36(1 - \xi)$  on Figure 7 (left).

<sup>d</sup> We will additionally require that the policy-maker (when imposing mandatory or optional XAI) wants to ensure that both firms enter the market. This last requirement helps because the binary choice of price may not be realistic in a monopoly.

Table 6. Equilibrium XAI levels in the simplified model under regulated and unregulated XAI

Horizontal Differentiation	Regulated XAI		Unregulated XAI
	Optional	Mandatory	
	$(\xi_1^*, \xi_2^*) = (\xi_{\text{policy}}^{\text{opt.}}, 0)$ where $\xi_{\text{policy}}^{\text{opt.}} = \min\{1, 2 - \eta\}$	$\xi_{\text{policy}}^{\text{mand.}} = \min\{1, 1 - \eta/2\}$	$(\xi_1^*, \xi_2^*) = \min\{1, 1 - \eta/2\}$
	$(\xi_1^*, \xi_2^*) = (0, \xi_{\text{policy}}^{\text{opt.}})$ where $\xi_{\text{policy}}^{\text{opt.}} = \delta$	No feasible $\xi_{\text{policy}}^{\text{mand.}}$	$(\xi_1^*, \xi_2^*) = (1 - \delta, 1)$

where  $\stackrel{\text{def}}{=} 1/2\beta t \times p/q$  and  $\delta \stackrel{\text{def}}{=} 4q/t \times (1 - \beta q^2/p)$ . An interesting finding is that **with or without regulations, the XAI market may receive less than full explanations** ( $\xi < 1$ ). When both mandatory and unregulated XAI are possible (the ‘max’ case), one might argue in favor of mandatory XAI as it ensures vertical XAI fairness by enforcing firms to offer an identical XAI level, even if it is not full explanation. That said, the results above show that **unregulated XAI can also deliver symmetric and therefore fair XAI**. In fact, the unregulated setting in ‘max’ is exactly the same as mandatory XAI—symmetric in price, quality, and XAI level. This is consistent with the indicative results from the full model discussed above. But contrary to the ‘max’ case, we find that no symmetric equilibrium is sustainable when firms use the same XAI methods (‘min’).

Overall, the results of our simplified model indicate that compared to regulated XAI, unregulated XAI may in fact provide better total welfare, total consumer utility, XAI fairness, and average XAI level<sup>a</sup>. While we cannot analytically make the unregulated and regulated comparison for the full model, the results with the simple model are unfavorable to a regulated setting because firms are unable to differentiate on price and quality to the same extent. Therefore, the results from the simple model are more relevant to markets where price and quality are relatively standardized, e.g., AI-based lending where the interest rates may be pegged across the market<sup>b</sup>.

<sup>a</sup> See Appendix A10.3.

<sup>b</sup> See Appendix A10.

## 7. Conclusion

Our paper provides several insights for policy-makers and managers. First, policy-makers in the field of XAI should pay attention to the differentiation of firms in terms of XAI methods and the market structure (explanation- vs. quality dominated). Firms may develop their own XAI algorithms in-house or customize existing packages to suit their needs. Whether firms differentiate or not in terms of the XAI method, warrants a different regulatory treatment. Another aspect of policy-making is the market structure that is formed as a result of the policy and firms' decisions. This study shows that in some cases, both markets are possible and firms have to choose one of them. Policy-makers should take this fact into account and carefully choose the required XAI level such that firms pick the market that yields higher total welfare. Sometimes under optional XAI policies, firms are indecisive about the market that they want to form and policy-makers should choose the required XAI level *low enough* that both firms choose the same market. In fact, one of our results is that high levels of explanations even under optional XAI may actually lead to no Nash equilibrium in the market. In addition, policy-makers must always pay attention to the marginal cost of quality and the cost of misfit explanations in the market, in order to find the optimal level of explanations.

The results confirm the intuition that mandatory XAI generally makes the society better off, and this holds regardless of firms' differentiation. That being said, in certain situations, **there is no additional benefit from mandating XAI**, especially when firms use similar XAI methods. This implies that policy-makers may need to focus on standardizing the level of explanations (optional XAI) or leaving the market unregulated, instead of enforcing firms to offer XAI.

Policy-makers should also notice that **requiring full explanations may actually make both firms and consumers worse off**, and this result holds for both optional and mandatory XAI policies. In fact, this is more of an issue under mandatory XAI. In this situation, policy-makers may have to settle for *partial* explanations, which means that AI models will remain partly blackbox. By contrast, optional XAI in differentiated markets may permit full explanations, albeit only one firm offering XAI. Therefore, we show that **there exists a tradeoff between maximizing the total welfare of the society (through mandatory XAI) and requiring full explanations (through optional XAI)**. We believe that this tradeoff is particularly important in differentiated markets because firms often develop their own XAI methods which are different from the competition (Bhatt et al., 2019).

An interesting observation in our results is that equilibrium solutions turn out to be expressed with  $\beta t$  occurring together. The interchangeability of  $\beta$  and  $t$  highlights the underlying tradeoff between quality and explanations and the fact that these two must be considered jointly in XAI decisions. To support this notion, notice that the low-quality firm in our model never benefits from offering XAI unilaterally. In fact, even in market E it is the *high-quality* firm that may benefit from XAI. Ironically, this finding points to the **significance of quality as the main product attribute for firms that want to offer XAI**. Thus, managers using ML algorithms in their products should focus their attention on quality and treat it as a necessary condition for offering profitable XAI.

Finally, in the absence of policy-makers, our analysis shows that when **firms are differentiated in XAI method, they should *mirror* each other’s explanation levels, but not necessarily offer full explanations**. It may appear that if technical threats of stolen IP, copycats, or adversarial attacks are eliminated and XAI methods are readily available for any AI algorithm, firms should maximize explanations as a lever to compete and charge higher prices. But our results show that even if these technical limitations are overcome, full explanations may not be the best strategy. The optimal explanation level depends on the misfit cost. As the misfit cost increases and the market cares more about explanations, managers should increase their XAI level by revealing more information but always hiding the rest.

## 8. Suggestions for Future Work

In this paper, we model XAI through an economics lens. Specifically, explanations are considered an additional component of AI-based products that increase customer utility. Future researchers should dig deeper than XAI level, and look into exact settings of XAI methods such as the number of features, R-squared explained, etc. Further, we do not model firms’ concerns about the confidentiality of their AI models. This factor may make firms even more reluctant to provide explanations. The accumulated knowledge in the model might be considered confidential and part of the firm’s algorithm property and trade secrets (Adadi and Berrada, 2018). The information revealed by XAI methods may be used to “steal” the underlying models (Barredo Arrieta et al., 2020; Orekondy et al., 2019) or launch adversarial attacks aimed at confusing the model (Goodfellow et al. 2015). For reasons such as these, firms may be less inclined to offer XAI in order to protect the privacy of their model. We believe further research is needed to investigate XAI in the presence of concerns about model privacy, adversarial attacks, and gaming by agents.



## 9. References

- Abdollahi, B., Nasraoui, O., 2016. Explainable Restricted Boltzmann Machines for Collaborative Filtering. ArXiv160607129 Cs Stat.
- Adadi, A., Berrada, M., 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). IEEE Access 6, 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- Ahmad, M.A., Eckert, C., Teredesai, A., McKelvey, G., 2018. Interpretable Machine Learning in Healthcare 7.
- Assad, S., Clark, R., Ershov, D., Xu, L., 2020. Algorithmic Pricing and Competition: Empirical Evidence from the German Retail Gasoline Market. SSRN Electron. J. <https://doi.org/10.2139/ssrn.3682021>
- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., Herrera, F., 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. Inf. Fusion 58, 82–115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- Bertini, M., Koenigsberg, O., 2021. The Pitfalls of Pricing Algorithms. Harv. Bus. Rev.
- Bhatt, U., Xiang, A., Sharma, S., Weller, A., Taly, A., Jia, Y., Ghosh, J., Puri, R., Moura, J.M.F., Eckersley, P., 2019. Explainable Machine Learning in Deployment.
- Bojarski, M., Del Testa, D., Dworakowski, D., Firner, B., Flepp, B., Goyal, P., Jackel, L.D., Monfort, M., Muller, U., Zhang, J., Zhang, X., Zhao, J., Zieba, K., 2016. End to End Learning for Self-Driving Cars. ArXiv160407316 Cs.
- Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., Elhadad, N., 2015. Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission, in: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '15. Association for Computing Machinery, New York, NY, USA, pp. 1721–1730. <https://doi.org/10.1145/2783258.2788613>
- Castelvecchi, D., 2016. Can we open the black box of AI? Nature 538, 20–23. <https://doi.org/10.1038/538020a>
- Che, Z., Purushotham, S., Khemani, R., Liu, Y., 2016. Interpretable Deep Models for ICU Outcome Prediction. AMIA Annu. Symp. Proc. AMIA Symp. 2016, 371–380.

- Csiszár, O., Csiszár, G., Dombi, J., 2020. Interpretable neural networks based on continuous-valued logic and multicriteria decision operators. *Knowl.-Based Syst.* 199, 105972. <https://doi.org/10.1016/j.knosys.2020.105972>
- Cunha, M., Osório C., A.M., Ribeiro, R.M., 2020. Endogenous Product Design and Quality When Consumers Have Heterogeneous Limited Attention (SSRN Scholarly Paper No. 2860456). Social Science Research Network, Rochester, NY. <https://doi.org/10.2139/ssrn.2860456>
- Dentons [WWW Document], 2021. URL <https://www.dentons.com/en/insights/guides-reports-and-whitepapers/2021/january/28/global-guide-to-autonomous-vehicles-2021> (accessed 12.16.21).
- Doshi-Velez, F., Kim, B., 2017. Towards A Rigorous Science of Interpretable Machine Learning. *ArXiv170208608 Cs Stat.*
- Doshi-Velez, F., Kortz, M., Budish, R., Bavitz, C., Gershman, S., O’Brien, D., Scott, K., Schieber, S., Waldo, J., Weinberger, D., Weller, A., Wood, A., 2017. Accountability of AI Under the Law: The Role of Explanation. <https://doi.org/10.48550/arXiv.1711.01134>
- Fu, R., Jin, G.Z., Liu, M., 2022. Human-Algorithm Interactions: Evidence from Zillow.com. Working Paper Series. <https://doi.org/10.3386/w29880>
- Goodfellow, I.J., Shlens, J., Szegedy, C., 2015. Explaining and Harnessing Adversarial Examples. *ArXiv14126572 Cs Stat.*
- Goodman, B., Flaxman, S., 2017. European Union Regulations on Algorithmic Decision-Making and a “Right to Explanation.” *AI Mag.* 38, 50–57. <https://doi.org/10.1609/aimag.v38i3.2741>
- Haspiel, J., Du, N., Meyerson, J., Robert Jr., L.P., Tilbury, D., Yang, X.J., Pradhan, A.K., 2018. Explanations and Expectations: Trust Building in Automated Vehicles, in: Companion of the 2018 ACM/IEEE International Conference on Human-Robot Interaction, HRI ’18. Association for Computing Machinery, New York, NY, USA, pp. 119–120. <https://doi.org/10.1145/3173386.3177057>
- Holzinger, A., Biemann, C., Pattichis, C.S., Kell, D.B., 2017a. What do we need to build explainable AI systems for the medical domain? *ArXiv171209923 Cs Stat.*

- Holzinger, A., Plass, M., Holzinger, K., Crisan, G.C., Pintea, C.-M., Palade, V., 2017b. A glass-box interactive machine learning approach for solving NP-hard problems with the human-in-the-loop. ArXiv170801104 Cs Stat.
- Howard, A., Zhang, C., Horvitz, E., 2017. Addressing bias in machine learning algorithms: A pilot study on emotion recognition for intelligent systems, in: 2017 IEEE Workshop on Advanced Robotics and Its Social Impacts (ARSO). Presented at the 2017 IEEE Workshop on Advanced Robotics and its Social Impacts (ARSO), IEEE, Austin, TX, USA, pp. 1–7. <https://doi.org/10.1109/ARSO.2017.8025197>
- Katuwal, G.J., Chen, R., 2016. Machine Learning Model Interpretability for Precision Medicine. ArXiv161009045 Q-Bio.
- Law Library of Congress (U.S.), G.L.R.D., , 2019. Regulation of artificial intelligence in selected jurisdictions. [WWW Document]. URL <https://purl.fdlp.gov/GPO/gpo123733>
- Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R.J., Wasserman, L., 2017. Distribution-Free Predictive Inference For Regression. ArXiv160404173 Math Stat.
- Lipton, Z.C., 2017. The Mythos of Model Interpretability. ArXiv160603490 Cs Stat.
- Lundberg, S.M., Lee, S.-I., 2017. A Unified Approach to Interpreting Model Predictions, in: Advances in Neural Information Processing Systems. Curran Associates, Inc.
- Malik, N., 2020. Does Machine Learning Amplify Pricing Errors in Housing Market?: Economics of ML Feedback Loops. <https://doi.org/10.2139/ssrn.3694922>
- McEneney, M.F., Kaufmann, K.F., 2005. Implementing the FACT Act: Self-Executing Provisions. *Bus. Lawyer* 60, 737–747.
- Miller, T., 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artif. Intell.* 267, 1–38. <https://doi.org/10.1016/j.artint.2018.07.007>
- Mittelstadt, B., Russell, C., Wachter, S., 2019. Explaining Explanations in AI, in: Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT\* '19. Association for Computing Machinery, New York, NY, USA, pp. 279–288. <https://doi.org/10.1145/3287560.3287574>
- Moore, J., Swartout, W., 1988. Explanation in Expert Systems: A Survey 58.

- Murdoch, W.J., Singh, C., Kumbier, K., Abbasi-Asl, R., Yu, B., 2019. Interpretable machine learning: definitions, methods, and applications. *Proc. Natl. Acad. Sci.* 116, 22071–22080. <https://doi.org/10.1073/pnas.1900654116>
- Neven, D., Thisse, J.-F., 1989. On Quality and Variety Competition.
- Orekondy, T., Schiele, B., Fritz, M., 2019. Knockoff Nets: Stealing Functionality of Black-Box Models. Presented at the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4954–4963.
- Ribeiro, M.T., Singh, S., Guestrin, C., 2016. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. *ArXiv160204938 Cs Stat.*
- Rudin, C., 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* 1, 206–215. <https://doi.org/10.1038/s42256-019-0048-x>
- Stanton, N.A., Salmon, P.M., Walker, G.H., Stanton, M., 2019. Models and methods for collision analysis: A comparison study based on the Uber collision with a pedestrian. *Saf. Sci.* 120, 117–128. <https://doi.org/10.1016/j.ssci.2019.06.008>
- Tan, S., Caruana, R., Hooker, G., Lou, Y., 2018. Distill-and-Compare: Auditing Black-Box Models Using Transparent Model Distillation. *Proc. 2018 AAAIACM Conf. AI Ethics Soc.* 303–310. <https://doi.org/10.1145/3278721.3278725>
- van Lent, M., Fisher, W., Mancuso, M., 2004. An explainable artificial intelligence system for small-unit tactical behavior, in: *Proceedings of the 16th Conference on Innovative Applications of Artificial Intelligence, IAAI’04.* AAAI Press, San Jose, California, pp. 900–907.
- Vandenbosch, M.B., Weinberg, C.B., 1995. Product and Price Competition in a Two-Dimensional Vertical Differentiation Model. *Mark. Sci.* 14, 224–249.
- Wang, Q., Huang, Y., Jasin, S., Singh, P.V., 2020. Algorithmic Transparency with Strategic Users (SSRN Scholarly Paper No. ID 3652656). Social Science Research Network, Rochester, NY. <https://doi.org/10.2139/ssrn.3652656>
- Wattal, S., Telang, R., Mukhopadhyay, T., 2009. Information Personalization in a Two-Dimensional Product Differentiation Model. *J. Manag. Inf. Syst.* 26, 69–95. <https://doi.org/10.2753/MIS0742-1222260204>
- Yurtsever, E., Lambert, J., Carballo, A., Takeda, K., 2020. A Survey of Autonomous Driving: Common Practices and Emerging Technologies. *IEEE Access* 8, 58443–58469. <https://doi.org/10.1109/ACCESS.2020.2983149>

Zhu, J., Liapis, A., Risi, S., Bidarra, R., Youngblood, G.M., 2018. Explainable AI for Designers: A Human-Centered Perspective on Mixed-Initiative Co-Creation, in: 2018 IEEE Conference on Computational Intelligence and Games (CIG). Presented at the 2018 IEEE Conference on Computational Intelligence and Games (CIG), IEEE, Maastricht, pp. 1–8. <https://doi.org/10.1109/CIG.2018.8490433>

## Appendix: Additional Results and Discussions

### A1. Proof of Lemma 1: Existence of Two Markets

Suppose that in addition to markets E and Q, a third market R exists, too. R can only be a combination of E and Q as follows:

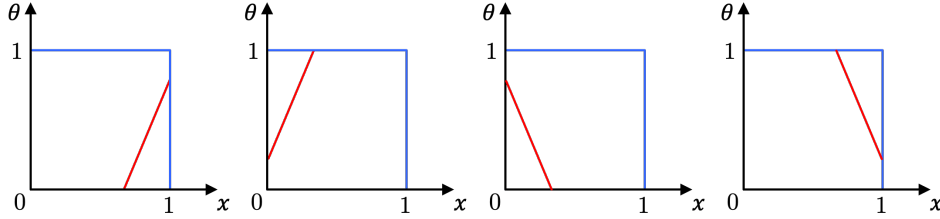
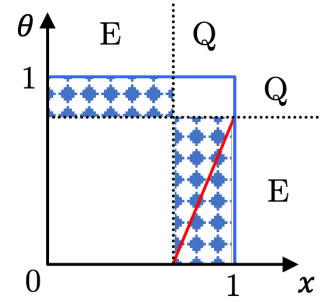


Figure 12. 4 possible combinations of markets E and Q

Now we show that none of these markets has a Nash equilibrium in our model. We run the proof for the first one; the rest are similar. Suppose that the first market in Figure 12 has a Nash equilibrium. Since preferences for quality and explanation are distributed uniformly throughout the unit square, one can break the hypothetical market R into submarkets E and Q<sup>a</sup>. The problem arises when we notice that two segments of R (shown by solid diamond grid) can be market E and Q at the same time, which is contradiction. In other words, no equilibrium exists for these regions and market R does not exist. ■



### A2. Backward Induction to Solve the Games

First, we calculate the equilibrium prices  $p_1^*$  and  $p_2^*$  in stage 2 by solving  $\partial\pi_1/\partial p_1 = 0$  and  $\partial\pi_2/\partial p_2 = 0$ , respectively. Next, we plug in these equilibrium prices in the profit functions and solve  $\partial\pi_1/\partial q_1 = 0$  and  $\partial\pi_2/\partial q_2 = 0$  to find the equilibrium quality levels  $q_1^*$  and  $q_2^*$ , respectively. Finally, we compare the profit functions (if available; otherwise, we use their partial derivatives with respect to  $\xi$ ) to determine which firm (if any) offers explanations. Doing this, we solve for the equilibrium decisions  $(s_1^*, s_2^*)$  in stage 0.

<sup>a</sup> Notice that breaking markets E and Q results in submarkets that are still E and Q, respectively.

### A3. Deriving Equilibria

We derive all equilibria twice: once by hand and once by writing a computer code that verifies our results. For the latter, we use the Python programming language along with a Computer Algebra System (CAS) called Sympy<sup>a</sup>, which is an open-source Python library for symbolic computation. Our Python code also enables us to run Monte-Carlo simulations when analytical solutions are not available (see Section 4.1.3).

### A4. Equilibrium Conditions

The conditions for equilibrium results mentioned in Table 2 are as follows:

Table 7. Equilibrium results mentioned in Table 2 along with their conditions

Firms' Horizontal Differentiation	
max	min
maxE1 Condition: $\beta t > \frac{1}{6(2\xi^2 - 7\xi + 6)}$ and $\xi > \xi_*$ $\beta t < 1/36(1 - \xi)$	minE1 Condition: $\beta t > 2/9\xi$
maxE0 Condition: $\beta t > \frac{1}{6(2\xi^2 - 7\xi + 6)}$ and $\xi \leq \xi_*$ $\beta t \geq 1/36(1 - \xi)$	minQ2 Condition: $\beta t < 4/9\xi$
maxQ2 Condition: $\beta t < 4/9\xi$	

### A5. Firms' Equilibrium Profits in Market E

Profit functions in market E are shown below. As we can see, they are somewhat more complicated than profits in market Q (See Table 3).

---

<sup>a</sup> <https://www.sympy.org/en/index.html>

Table 8. Firms' equilibrium profits in market E

Firms' Horizontal Differentiation		$\pi_1$	$\pi_2$
	max	If $\xi > \xi_*$ : $\frac{(12\beta t\xi - 36\beta t + 1)^2(1 - 36\beta t\xi + 72\beta t)}{144\beta(18\beta t\xi - 36\beta t + 1)^2}$	If $\xi \geq \xi_*$ : $\frac{(24\beta t\xi - 36\beta t + 1)^2(1 - 36\beta t\xi + 72\beta t)}{144\beta(18\beta t\xi - 36\beta t + 1)^2}$
		If $\xi < \xi_*$ : $\frac{t}{2} - \frac{1}{144\beta}$	If $\xi < \xi_*$ : $\frac{t}{2} - \frac{1}{144\beta}$
	min	$\frac{(18\beta t\xi - 1)(12\beta t\xi - 1)^2}{144\beta(9\beta t\xi - 1)^2}$	$\frac{(18\beta t\xi - 1)(6\beta t\xi - 1)^2}{144\beta(9\beta t\xi - 1)^2}$

## A6. Proof of Lemma 2

First, notice that both minE1 and minQ2 are possible when  $2/9 < \beta t\xi < 4/9$ . We start by proving that in this region, minE1 profits are always increasing in  $\xi$  for both firms. Then, given that profits are constant in minQ2, we show that  $\pi_1^{\min E1}(\beta t\xi = 2/9) > \pi_1^{\min Q2}(\beta t\xi = 2/9)$ , indicating that firm 1 always earns more in minE1 than in minQ2. Conversely, we show that  $\pi_2^{\min E1}(\beta t\xi = 4/9) < \pi_2^{\min Q2}(\beta t\xi = 4/9)$ , implying that firm 2 always earns less in minE1 than in minQ2.

### A6.1. Part 1 of the Proof

Taking the derivative of  $\pi_1$  w.r.t.  $\xi$  we get:

$$\frac{\partial \pi_1}{\partial \xi} = \frac{t(12\psi - 1)(324\psi^2 - 81\psi + 4)}{24(9\psi - 1)^3}$$

where  $\psi := \beta t\xi$ . The terms  $(12\psi - 1)$  and  $(9\psi - 1)$  are positive for  $\psi > 2/9$ . Next, we find the roots of  $324\psi^2 - 81\psi + 4 = 0$  and find that this expression is positive after its second root  $\psi_2 = 0.81 < 2/9$ . Therefore, the whole derivative is positive for the given region of  $\beta t\xi$ . Similar result holds for  $\partial \pi_2 / \partial \xi$ :

$$\frac{\partial \pi_2}{\partial \xi} = \frac{t(6\psi - 1)(162\psi^2 - 27\psi + 2)}{24(9\psi - 1)^3}$$

which is always positive in the given range of  $\beta t\xi$ .



## A6.2. Part 2 of the Proof

Remember that  $\pi_1^{\min E1} = (12\psi - 1)^2(18\psi - 1)/(144\beta(9\psi - 1)^2)$ . At  $\psi = 2/9$  we get  $\pi_1^{\min E1}(\psi = 2/9) \approx 0.06/\beta$ , which is already greater than  $\pi_1^{\min Q2} = 4/81\beta \approx 0.05/\beta$ .

As with the second firm,  $\pi_2^{\min E1} = (6\psi - 1)^2(18\psi - 1)/(144\beta(9\psi - 1)^2)$ . At  $\psi = 4/9$  we have  $\pi_2^{\min E1} \approx 0.015/\beta$ , which is still less than  $\pi_2^{\min Q2} = 2/81\beta \approx 0.025/\beta$ . ■

## A7. Analytical Expressions of Total Welfares

Table 9. Analytical expression of total welfare at each equilibrium

Equilibrium	Total Welfare ( $W_t$ )
maxE0	$\frac{1 - 18\beta t}{36\beta}$
maxE2	$\frac{9\beta t(\xi - 1) + 1}{36\beta}$
maxE1	Too Long; available in our computer code.
maxQ2	$\frac{27\beta t(\xi - 1) + 2}{54\beta}$
minQ2	$\frac{3.375e32\beta t\xi - 3.375e32\beta t + 5.0e31}{1.35e33\beta}$ where $xy$ means $x \times 10^y$ .
minE1	$\frac{1296\beta t^3\xi^3 - 1458\beta t^3\xi^2 - 90\beta t^2\xi^2 + 324\beta t^2\xi - 21\beta t\xi - 18\beta t + 2}{72\beta(9\beta t\xi - 1)^2}$

## A8. Analysis of Mandatory XAI Policies

The **red** region ① in Figure 8 (left) is where firms earn negative profits in maxE2. This is found by setting  $\pi_1 = \pi_2 = t(1 - \xi)/2 - 1/144 < 0$  and solving for  $\beta t$ . The expression for the orange region can be found by comparing the smallest profit that firms can obtain in maxQ2 with its corresponding profit in maxE2. In this case, we need to compare  $\pi_2^{\max Q2} = 2/81\beta$  with  $\pi_2^{\max E2} = t(1 - \xi)/2 - 1/144$  and solve for  $\beta t$ .

In the **purple** area, only market Q is possible. Since  $W_t$  is positively correlated with  $\xi$ , the policy-maker chooses the maximum possible  $\xi$  for any given  $\beta t$ . Therefore, for  $\beta t < 1/6(2 \cdot 1^2 - 7 \cdot 1 + 6) = 1/6$  the policy-maker sets  $\xi_{\text{policy}}^{\text{mand.}} = 1$  so that the market becomes Q and maximum  $W_t$  is obtained. If the policy-maker would choose smaller  $\xi$  (e.g.,  $\xi = 0.5$ ), then for some values of  $\beta t$ , firms

would choose market E to earn higher profits and maximum total welfare would not be achieved.

Inside the **green** and **orange** regions, both markets E and Q are possible, but firms choose market Q in the orange region and E in the green region to maximize their profits. For  $1/6 = A \leq \beta t < B = 4/9$ , the policy-maker still sets  $\xi_{\text{policy}}^{\text{mand.}} = 1$  so that the market remains Q and maximum  $W_t$  is achieved. From point B to C, the policy-maker must choose  $\xi_{\text{policy}}^{\text{mand.}} = 4/9\beta t - \varepsilon$  where  $\varepsilon \rightarrow 0$ , the reason being that  $\xi = 4/9\beta t$  borders on the red region ① in which  $\pi_1, \pi_2 < 0$  and firms leave the market, resulting in  $W_t = 0$ . Thus, with the above  $\xi_{\text{policy}}^{\text{mand.}}$ , the market remains Q from B to C and firms earn positive profits.

In the **blue** region, the policy-maker chooses  $\xi_{\text{policy}}^{\text{mand.}} = 1 - (1/72)/\beta t$  from point C to D onward and the market is now E. Under this policy, from point C onward, firms earn zero profit and they cannot form a Q market from C to D due to the choice of  $\xi_{\text{policy}}^{\text{mand.}}$  by the policy-maker.

Table 10. 5 regions on the  $\beta t$ - $\xi$  plane under mandatory XAI and firms use different XAI methods

Region(s)	Which Market is Possible?	Relation Between $\beta t$ and $\xi$
Blue and Red	only E	In Red and Blue: $\beta t > 4/9\xi$
		In Red: $\beta t < (1/72)/(1 - \xi)$
Green and Orange	E and Q	In Green and Orange: $1/6(2\xi^2 - 7\xi + 6) < \beta t < 4/9\xi$
		In Orange: $\beta t < 0.063/(1 - \xi)$
Purple	only Q	$\beta t < 1/6(2\xi^2 - 7\xi + 6)$

## A9. Number of Firms that Offer XAI w.r.t. $\xi$

Figure 11 is obtained from Figure 7 as follows: When firms use the same XAI method, for small  $\beta t$  both firms offer XAI regardless of  $\xi$ —the green region in Figure 7, right. Therefore, the number of firms that offer XAI is constant in this case. Now let us increase  $\beta t$ . If  $\xi$  is small, then both firms offer XAI, but if it is

large, then only one firm offers XAI. In the gray region, there is no equilibrium, and no firm offers XAI.

On the other hand, when firms use different XAI methods, Figure 7, left, tells us that with small  $\beta t$ , both firms offer XAI for all values of  $\xi$ —the green region in the figure. As  $\beta t$  increases, there is a region where increasing  $\xi$  results in switching from “No Nash Equilibrium” (gray region) to both firms offering XAI. After that, increasing  $\beta t$  results in a situation where no firm offers XAI (the blue and gray regions in Figure 7) for all  $\xi$ . Finally, as  $\beta t$  keeps increasing, we see that at most one firm offers XAI—the yellow region.

## A10. Unregulated XAI Level

### A10.1. maxE (Full Model)

Firms play the same game as in the exogenous case, except that their decisions in stage-0 are not about offering/not offering XAI, but rather about the level of explanations that they offer. We can solve this by plugging  $p^*$  and  $q^*$  in the profit functions and solving  $\partial\pi_1/\partial\xi_1 = 0$  and  $\partial\pi_2/\partial\xi_2 = 0$  to find the stage-0 equilibrium levels of explanations  $\xi_1^*$  and  $\xi_2^*$ , respectively. The equilibrium results are:

$$\begin{aligned}
p_1^* &= \frac{-t(\xi_1 + \xi_2 - 2)(12\beta t(\xi_1 + 2\xi_2 - 3) + 1)}{2(18\beta t(\xi_1 + \xi_2 - 2) + 1)}; \\
p_2^* &= \frac{-t(\xi_1 + \xi_2 - 2)(12\beta t(\xi_2 + 2\xi_1 - 3) + 1)}{2(18\beta t(\xi_1 + \xi_2 - 2) + 1)}; \\
q_1^* &= \frac{12\beta t(\xi_1 + 2\xi_2 - 3) + 1}{12\beta(18\beta t(\xi_1 + \xi_2 - 2) + 1)}; \\
q_2^* &= \frac{12\beta t(\xi_2 + 2\xi_1 - 3) + 1}{12\beta(18\beta t(\xi_1 + \xi_2 - 2) + 1)}; \\
\xi_1^* &= 1 - \frac{1}{36\beta t}; \\
\xi_2^* &= 1 - \frac{1}{36\beta t}; \\
\Delta p &= \frac{6\beta t^2(\xi_1 - \xi_2)(\xi_1 + \xi_2 - 2)}{18\beta t(\xi_1 + \xi_2 - 2) + 1}; \\
\Delta q &= \frac{-t(\xi_1 - \xi_2)}{18\beta t(\xi_1 + \xi_2 - 2) + 1}; \\
\pi_1 &= -\frac{(12t\beta\xi_1 + 24t\beta\xi_2 - 36t\beta + 1)^2(36t\beta\xi_1 + 36t\beta\xi_2 - 72t\beta + 1)}{144\beta(18t\beta\xi_1 + 18t\beta\xi_2 - 36t\beta + 1)^2}; \\
\pi_2 &= -\frac{(24t\beta\xi_1 + 12t\beta\xi_2 - 36t\beta + 1)^2(36t\beta\xi_1 + 36t\beta\xi_2 - 72t\beta + 1)}{144\beta(18t\beta\xi_1 + 18t\beta\xi_2 - 36t\beta + 1)^2}.
\end{aligned}$$

Notice that this is a symmetric equilibrium in  $p^*$ ,  $q^*$ , and  $\xi^*$ .

### A10.2. minE (Full Model)

If we solve for equilibrium XAI levels in this case, we obtain the following system of equations:

$$\begin{aligned}\xi_1 &= \frac{12\beta t \xi_2 + 1}{12\beta t}; \\ \xi_2 &= \frac{6\beta t \xi_1 - 1}{6\beta t} \Rightarrow \xi_1^* = \xi_2^* + \frac{1}{6\beta t}\end{aligned}$$

which is overdetermined and has no solution. Therefore, sustaining market E while using the same XAI method is not possible.

### A10.3. Simplified Model

The game still has three stages:

- Stage 0: firms choose  $\xi_1^*, \xi_2^* \in [0, 1]$ ;
- Stage 1: firms choose  $q_1^*, q_2^* \in \{0, q\}$ ;
- Stage 2: firms choose  $p_1^*, p_2^* \in \{0, p\}$ .

We only examine equilibria where both firms enter the market, i.e., they choose  $p_1^* = p_2^* = p$  in stage 2. Again, backward induction is used to solve for equilibrium qualities and XAI levels. In each stage, we find firm  $i$ 's best response to the other firm's action.

#### A10.3.1. Firms Differentiated ('max) in XAI Methods

Only the explanation-dominated market is viable. We identify an equilibrium where unregulated firms offer partial XAI ( $\xi < 1$ ) in equilibrium:

$$p_1^* = p_2^* = p; \quad q_1^* = q_2^* = 0; \quad \xi_1^* = \xi_2^* = \min\{1, 1 - \eta/2\}$$

where  $\eta \stackrel{\text{def}}{=} 1/2\beta t \times p/q$ . This equilibrium corresponds to maxE2 in the full model and is sustained if  $p/\beta q^2 \in [1, 2]$  and  $\beta q \leq 1/4$ . The total welfare  $W_t$  and total consumer utility in unregulated XAI are better than regulated XAI when  $\eta \in [1.33, 2]$ , worse when  $\eta \in [1, 1.33]$ , and the same otherwise. To get a sense of this equilibrium, see firm payoffs in Table 11.

Table 11. The game between firms in stage 2 of the simplified model. In each cell, the top row is  $\pi_1$  and the bottom row  $\pi_2$ .

		Firm 1's quality choice ( $q_1$ )	
		0	$q$
Firm 2's quality choice ( $q_2$ )	0	$\xi_r p$ $(1 - \xi_r) p$	$(\xi_r + q/2t\bar{\xi})p - \beta q^2$ $\left((1 - \xi_r) - q/2t\bar{\xi}\right)p$
	$q$	$(\xi_r - q/2t\bar{\xi})p$ $\left((1 - \xi_r) + q/2t\bar{\xi}\right)p - \beta q^2$	$\xi_r p - \beta q^2$ $(1 - \xi_r)p - \beta q^2$

where  $\bar{\xi} \stackrel{\text{def}}{=} 1 - \xi_1 + 1 - \xi_2$  and  $\xi_r \stackrel{\text{def}}{=} (1 - \xi_2)/\bar{\xi}$ . With symmetric XAI ( $\xi_1 = \xi_2, \xi_r = 0.5$ ) and zero quality ( $q_1 = q_2 = 0$ ), each firm earns a payoff of  $p/2$ . As firm 1 (firm 2) increases XAI  $\xi_r$  increases (decreases) and so does firm 1's (firm 2's) market share. If firm 1 increases quality ( $q_1 = q$ ) they get additional market share  $q/2t\bar{\xi}$  but pay cost  $\beta q^2$ .

### A10.3.2. Firms Undifferentiated ('min') in XAI Methods

Similar setup as before. Only quality dominated Q-market is viable. We identify an equilibrium where unregulated firms offer less than full XAI in equilibrium:

$$p_1^* = p_2^* = p; \quad (q_1^*, q_2^*) = (q, 0); \quad (\xi_1^*, \xi_2^*) = (1 - \delta, 1)$$

where  $\delta \stackrel{\text{def}}{=} 4q/t \times (1 - \beta q^2/p)$ . This equilibrium corresponds to minQ2 in the full model and holds when  $p/\beta q^2 \in [1, 2]$ . Since mandatory XAI is not feasible, we only compare unregulated and optional XAI. Interestingly, total welfare, total consumer utility, and average XAI level are all higher in unregulated XAI than in optional XAI.

---

## Web Links

<sup>1</sup> <https://www.gartner.com/en/information-technology/insights/top-technology-trends>

<sup>2</sup> [https://www.idc.com/getdoc.jsp?containerId=IDC\\_P33198](https://www.idc.com/getdoc.jsp?containerId=IDC_P33198)

<sup>3</sup> <https://www.newyorker.com/magazine/2017/04/03/ai-versus-md>

<sup>4</sup> <https://www.darpa.mil/program/explainable-artificial-intelligence>

<sup>5</sup> <https://ec.europa.eu/futurium/en/ai-alliance-consultation.1.html>

- 
- <sup>6</sup> <https://tinyurl.com/theguardian-AI-watchdog>
- <sup>7</sup> <https://research.aimultiple.com/xai/>
- <sup>8</sup> <https://tinyurl.com/theguardian-transparent-ai>
- <sup>9</sup> <http://www.ftc.gov/legal-library/browse/statutes/fair-credit-reporting-act>
- <sup>10</sup> <https://insight.equifax.com/visualizing-xai-in-credit-risk/>
- <sup>11</sup> <https://www.nts.gov/news/events/Pages/2020-HWY18FH011-BMG.aspx>
- <sup>12</sup> <https://www.compact.nl/en/articles/autonomous-compliance/>
- <sup>13</sup> <https://www.wired.com/story/europes-new-privacy-law-will-change-the-web-and-more/>
- <sup>14</sup> <https://tinyurl.com/vox-GDPR-data-protection>
- <sup>15</sup> <https://tinyurl.com/cnet-GDPR-fb-google>
- <sup>16</sup> <https://tinyurl.com/forbes-GDPR-benefits>
- <sup>17</sup> <https://tinyurl.com/vox-GDPR-data-protection>
- <sup>18</sup> <https://tinyurl.com/cnbc-fb-call-for-GDPR>
- <sup>19</sup> <https://tinyurl.com/comply-or-explain>
- <sup>20</sup> [https://www.wipo.int/news/en/wipolex/2016/article\\_0014.html](https://www.wipo.int/news/en/wipolex/2016/article_0014.html)
- <sup>21</sup> <https://tinyurl.com/GDPR-harmful-for-AI>
- <sup>22</sup> <https://tinyurl.com/cnbc-Intel-AI-infancy>