# Incremental XAI: Memorable Understanding of AI with Incremental Explanations

Jessica Y. Bo
jbo@cs.toronto.edu
University of Toronto*
Toronto, Canada

Pan Hao
panhao@u.nus.edu
National University of Singapore
Singapore, Singapore

Brian Y. Lim
brianlim@comp.nus.edu.sg
National University of Singapore
Singapore, Singapore

## ABSTRACT

Many explainable AI (XAI) techniques strive for interpretability by providing concise salient information, such as sparse linear factors. However, users either only see inaccurate global explanations, or highly-varying local explanations. We propose to provide more detailed explanations by leveraging the human cognitive capacity to accumulate knowledge by incrementally receiving more details. Focusing on linear factor explanations (factors × values = outcome), we introduce Incremental XAI to automatically partition explanations for general and atypical instances by providing Base + Incremental factors to help users read and remember more faithful explanations. Memorability is improved by reusing base factors and reducing the number of factors shown in atypical cases. In modeling, formative, and summative user studies, we evaluated the faithfulness, memorability and understandability of Incremental XAI against baseline explanation methods. This work contributes towards more usable explanation that users can better ingrain to facilitate intuitive engagement with AI.

## CCS CONCEPTS

• **Human-centered computing → Empirical studies in HCI**; • **Computing methodologies → Artificial intelligence**.

## KEYWORDS

explanations, explainable AI, cognitive load, memory

## 1 INTRODUCTION

As Artificial Intelligence (AI) systems become prevalent, it is paramount for explainable AI (XAI) to be developed to support their proper use and understanding [1, 3, 4, 33, 38, 55]. Although much work has shown that XAI can improve satisfaction and trust [26,

*Work was performed while at the National University of Singapore

33, 41, 45, 60], many studies have failed to demonstrate measurably improved understanding [24, 48]. This requires users to ingrain AI explanations to quickly recall and apply the knowledge for decision making. Providing short explanations like sparse linear models could help, but these would be too simplified to be faithful to the complex underlying AI decision, and mislead users [45]. In contrast, more expressive explanations may be more faithful, but can be challenging to read or recall [2], hindering their accessibility. This is especially important for users to generalize their understanding of AI behavior for future scenarios [44]. Hence, XAI needs to be sufficiently detailed, yet memorable to support effective understanding.

To help users develop a richer understanding of AI models, instead of inundating users with complex explanations, we propose to explain *incrementally*. This is inspired from pedagogy, where students learn a concept gradually rather than all-at-once. For example, physics students learn about classical Newtonian mechanics for objects moving at common speeds, but later learn the theory of Special Relativity that describes objects at very high speeds with the Lorentz transformation. Understanding relativistic mechanics is very complicated and requires the foundational understanding of classical mechanics first. Thus, we argue that users can eventually understand complex explanations and models, but they should be grounded on simpler explanations, and incrementally informed.

We propose a step towards elevating user understanding towards complex AI explanations with *Incremental XAI*. This framework defines how to explain AI predictions from typical cases to outlier cases. We investigate this for simple surrogate explanation models, specifically, sparse linear explanations that describe linear factors that multiply against feature values. For example, a factor of $w^{(\text{Bathrooms})} = \$17k$ explains the predicted price of a house based on # Bathrooms by indicating that each bathroom adds $17k, and that two bathrooms would contribute $34k together. We begin with partitioning the dataset into subspaces, training a linear model in the majority (typical) subspace with Base factors, and training a linear model in the minority (outlier) subspace with (Base + Incremental) factors. To minimize new information to learn, we regularize the Incremental factors to be 0 when possible. In our bathrooms example, the majority smaller houses could have a rate of $17k/bathroom, while minority larger houses can have costlier bathrooms at $17k + $51k, perhaps due to luxury fittings. We contribute:

- The *Incremental XAI* paradigm which enables gradual delivery of complex explanations, gaining the benefit of multiple lightweight explanations that achieves higher faithfulness.
- A tree-based incremental explanation using linear model trees, additive factors, and factor sparsity regularization. We also developed a tabular user interface to convey explanations incrementally, and contrasted this with baseline variants.

- An evaluation of the faithfulness, usage, understanding, and memorability of Incremental explanations against Global, Subglobal, and Local baseline explanations in modeling, formative, and summative user studies. We compared Incremental explanations with Global explanations to evaluate if providing more detailed explanations based on category of cases (subspaces) helps understanding, with Subglobal explanations that are a baseline subspace model that explains each subspace independently, and with Local explanations since they are often singularly deployed primarily for instance-based explanations, but may be misused for general understanding.
- A discussion of how to generalize the Incremental XAI paradigm to other applications and AI explanations.

## 2 BACKGROUND AND RELATED WORK

Explainable AI (XAI) remains problematic for human interpretation due to the inaccuracy of overly simplified methods. Here, we give a primer on XAI and and their cognitive demands, techniques to mitigate cognitive load, the need to provide multiple explanations, and XAI techniques partitioned into subspaces to improve accuracy.

### 2.1 Surrogate explanations of AI

Explanations of AI can improve user understanding by providing surrogate explanations of accurate AI models, or making "glassbox" models that are intrinsically interpretable. However, the latter approach may have limited accuracy since these models tend to be overly simple. Instead, we focus on providing surrogate explanation models that approximate complex AI models that retains the use of accurate AI models while explaining with some unfaithfulness.

Miller [44] identified two goals for explanations in AI: i) to select a small set of causes for an observation [39], and ii) generalize observations into a conceptual model to predict and control future cases [22]. Wang et al. identified other reasoning processes that XAI should support [63]. Our research objective is to improve XAI techniques to better support the second goal of a generalized understanding. This requires explanations to be intuitive and memorable so that users can rapidly apply their knowledge to anticipate the AI model's behavior in future settings. Global explanations provide a suitable basis to support this goal. They answer the question "*How* does the AI model make predictions?" Techniques include explaining the AI decision in terms of linear factors [48], nonlinear partial dependence plots [27] and generalized additive models [2, 12], and decision trees [49] and rules [29, 32, 52].

To support the former goal of explaining causes for an individual case, instance explanations are provided instead. These answer the question: "*Why* did the AI model make this prediction?" Techniques include feature attributions [5, 15, 42, 59], and counterfactual explanations [11, 62]. An instance explanation only explains the decision for a target instance and may provide a different explanation for another instance; thus it may not generalize to multiple instances. To overcome this, Ribeiro et al. proposed local explanations to train explainer models on instances similar to the target instance of interest [51]. Since these explanations focus on narrower sets of instances, they are more faithful to the underlying AI being explained, but require users to remember many models for dissimilar instances.

Given their ubiquity in XAI practice, we include them in our investigations. In this work, we aim to provide explanations that are memorable like Global explanations and faithful like Local explanations, investigated Subglobal explanations that balance between the two, and proposed Incremental explanations that improve the memorability of Subglobal explanations.

### 2.2 Cognitive demands of AI Explanations

Although XAI aims to improve user understanding, they are not necessarily easy for users to interpret [34]. High cognitive load harms user experience and the effectiveness of AI explanations [28, 48]. This is often measured by the number of attributes used in explanations [16] or the nonlinearity of visualizations [2]. Indeed, people consider simpler explanations as more probable than those with more clauses [40], but oversimplifying explanations will erode trust in XAI [47, 67]. Explanations need to be delivered at the right level of cognitive effort to be effective [20, 25, 31, 58].

A simple method to get users to understand explanations is to prompt users to think when reading explanations [10], but this does not ensure deep learning and understanding or make explanations less cognitively demanding. Several techniques have been proposed to reduce cognitive load. The most common is to do feature reduction to limit the number of attributes shown to users. This can be accomplished with sparsity regularization [45] and constraining explanations to use integer coefficients instead of real numbers [61]. However, this limits the expressiveness of explanations that users could consume. Another approach is to simplify more sophisticated visual explanations, such as nonlinear line graphs. Cognitive-optimized GAM (COGAM) balances cognitive load and accuracy by quantifying the visual cognitive chunks in line chart explanations, and providing a hybrid explanation with sparse linear factors and less curvy line charts [2]. However, these approaches only optimize one explanation at a time, but neglects the human cognitive capacity to accumulate knowledge.

### 2.3 Providing multiple AI explanations

Accurate understanding of an AI system requires detailed knowledge of its parameters and non-linear decisions, yet explanations need to be simple for easy comprehended. To avoid information overload, detailed explanations can be provided on demand [36] or with progressive disclosure [58]. Users have various demands for explanations [35], diverse usage strategies of explanations [36], and use multifaceted explanations to understand AI decisions [37, 66].

Hence, instead of considering XAI interpretation as independent interactions, it should be considered as sequentially dependent accumulation of knowledge (e.g., dialogic [44]). People use explanations to build a mental model of the AI, so successful explanations can be measured through the goodness of the learned mental model [33, 65]. Mental models play a key role in human-AI interactions [6], but can be formed poorly without intervention [18]. In this work, we propose a new paradigm of providing explanations *incrementally* by ensuring that the shallower, simpler and explanations can smoothly transition into deeper, more detailed ones. This leverages the human ability for cumulative learning [57], and allows users to understand how the explanations relate [68] at different levels.
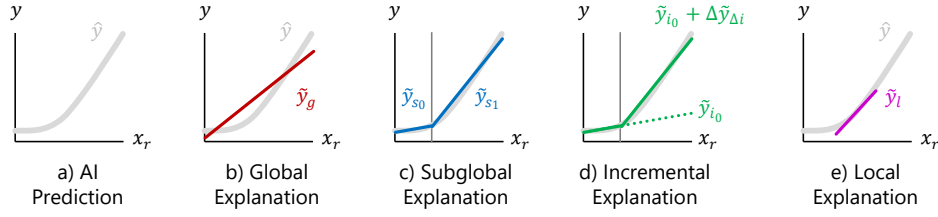
**Figure 1: Conceptual examples of XAI types with univariate (1D) data shown for simplicity; see Fig. A.1 for 2D multivariate examples with real data. a) Original AI System predicts output $\hat{y}$ non-linearly with respect to attribute $x_r$. b) Global explainer that approximates $\hat{y}$ with a linear equation $\tilde{y}_g \propto x_r$. c) Subglobal explainer increases faithfulness by segmenting along $x_r$ to provide multiple linear explanations $\tilde{y}_{s_1} \propto x_r, x_r < \chi_r$ and $\tilde{y}_{s_2} \propto x_r, x_r \geq \chi_r$. d) Incremental explainer that is similar to Subglobal, but first explains with a linear model $\tilde{y}_{i_0} \propto x_r$ the contiguous majority of instances (in this case, $x_r < \chi_r$), then explains outlier instances ($x_r \geq \chi_r$) with an additive linear model $\tilde{y}_{i_0} + \Delta\tilde{y}_{\Delta i}$. e) Local explanation explains each instance with a linear equation $\tilde{y}_l \propto x_r$ based on neighboring instances. Multiple local explanations are needed to represent the full input space.**

**Table 1: Comparison of linear explanation models with varying faithfulness and memorability due to the # factors used, which affects their expressiveness and the # terms for the user to remember. The AI System typically has too many factors to be interpretable, while sparse linear explanations consider a sparse set of $n$ factors. A Global explanation and Local explanation have 1 set of $n$ factors, but the latter requires many $N$ explanations to understand all use cases cumulatively. Subglobal explanations split the instances into $k$ subspaces, needing $kn$ factors to explain fully. Incremental explanations similarly can explain the same $k$ subspaces, but reuses some factors, and can omit negligible changes, thus it has $\leq kn$ factors.**

|  | Prediction Model | Explainer Models | | | |
|---|---|---|---|---|---|
|  |  | Global | Subglobal | **Incremental** | Local |
| # Factors | $\gg n$ | $n$ | $kn, k \geq 2$ | $\leq kn$ | $Nn, N \gg n$ |
| Faithfulness | Self | Low | Med | Med | High |
| Memorability | Low | High | Med-Low | Med-High | Low |

## 2.4 Subspace-based XAI techniques

Several XAI techniques have been developed to address the short-comings of global explanations beingf too coarse and local explanations being too narrow. We discuss methods that divide instances into subspaces and explain each subspace separately. Methods are based on trees, rules, or aggregation.

Model agnostic multilevel explanations (MAME) [46] provides an explanation tree with weights at each node, representing a progression from a global explanation at the root to local explanations at the leaves. However, their method does not enforce stability between each linear model, so it would be difficult for users to learn each sub-explanation incrementally. Model Understanding through Subspace Explanations (MUSE) [30] provide decision sets for different subspaces by simultaneously optimizing for faithfulness and rule compactness, but the explanations are in terms of rules unlike our factors-based format, and the attributes are not necessarily the same for each subspace, thus not consistent. Equi-explanation Maps [14] divide a feature space into hyper-cuboid subspaces (i.e., defined within min/max ranges for specific attributes) that are consistent, and explain each subspace with linear classifiers, but its boundary definitions are much more complex. Submodular Pick LIME [51] leverages instance-based local LIME explanations to provide a global explanation by picking diverse LIME explanations that have high non-redundant coverage. This aims to limit the total number of Local explanations needed to achieve global understanding. GLObal to loCAL eXplainer (GlocalX) [56] iteratively merges

local decision rules into global explanations to provide a smooth pathway from detailed local explanations to more general global explanations and vice versa. This was demonstrated for rule-based explanations of classification, unlike our prediction regression task with linear factors. Sparse LInear Subset Explanations (SLISE) [7] is a robust regression method that finds the largest subset in the data and trains a sparse linear model. Our method could use this to learn the base model of the Incremental explanation. SLISEMAP [8] extends SLISE to group instances into clusters based on the similarity of their local explanations. This involves dimensionality reduction, so the resulting dimensions are not explicitly interpretable.

All these methods aim to explain each subspace faithfully, but neglect to account for users having to remember or relate across subspaces. Thus the inter-subspace consistency is low. In our work, we focus on first providing a base explanation for a majority subspace, and explain remaining smaller subspaces as incrementally different from the base. This new requirement stems from usability needs of XAI that the prior works neglect.

## 3 TECHNICAL APPROACH

We first describe baseline explanation approaches using sparse linear factors, then articulate our Incremental explanation approach. An AI System's prediction $\hat{y}$ is typically generated from many input attributes (features) $\boldsymbol{x}$, and the AI's decision may change non-linearly with each attribute (e.g., price can increase exponentially with living area of a house). Sparse linear models provide simple
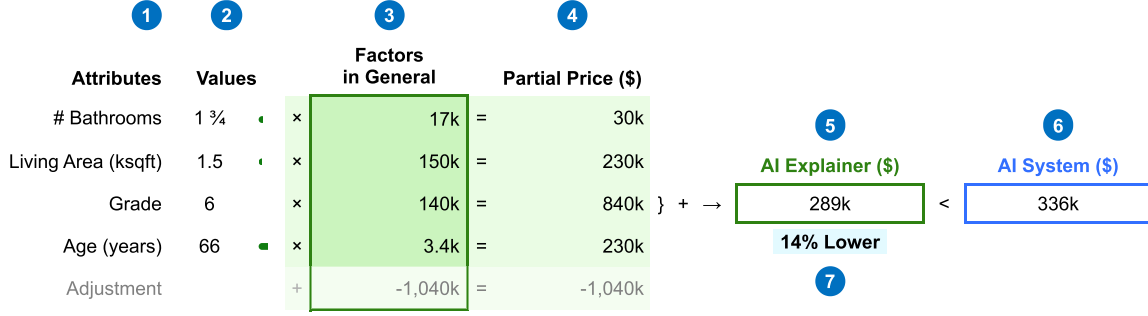
**Figure 2: User interface (UI) of AI System with Global explanation showing: 1) attributes used for prediction, 2) their values $x^{(r)}$ for the given instance, 3) factors $w_r$ that the explainer multiplies with values, 4) partial contributions $\tilde{y}_i = w^{(r)}x^{(r)}$ of each attribute, 5) output estimation $\tilde{y} = \sum_r \tilde{y}^{(r)}$ from the AI Explainer, 6) prediction $\hat{y}$ from the AI System, with inequality indication ($<$ in this case), and 7) indicator of how different the AI Explainer estimation is from the AI System prediction. Factors are the same for all instances and do no change. Different information may be hidden under various test conditions.**

explanations by articulating only a few important attributes (hence sparse), and indicating how each attribute influences the prediction. These are typically presented as *feature attributions*, i.e., positive or negative numbers indicating the direction and magnitude of the influence. However, attributions are not particularly easy to track or interpret, since they vary inconsistently for different instances.

Instead, like Poursabzi-Sangdeh et al. [48], we focus on sparse linear factor explanations that compute the feature attribution of the $r$th attribute $y^{(r)}$ as a multiplication of a factor weight $w^{(r)}$ and the attribute value $x^{(r)}$, i.e., $\tilde{y}^{(r)} = w^{(r)}x^{(r)}$. For example, consider a house with 1¾ bathrooms, i.e., $x^{(1)} = 1.75$. $w^{(1)} = 17$k means that each increase in one bathroom costs $17k more, so the # Bathroom contributes $17k × 1.75 = $30k to the total house price. Users can apply these factors to another instances to calculate how the AI Explainer would estimate the AI prediction for that instance. For example, a house with 3 bathrooms would have its # Bathrooms contribute $17k × 3 = $51k to its price. Sparse linear factors can be applied broadly to all instances (Global explanation), semi-broadly to groups of instances (Subglobal explanation), or to individual instances (Local explanation). We introduce each of these explanation methods and then describe our approach for the Incremental explanation that extends Subglobal. Each type has varying *faithfulness* to the AI prediction and *memorability* for users to recall the factors, which we summarize in Table 1 and illustrate conceptually in a 1-dimensional example in Fig. 1. We refer to these explanation variants as XAI types.

## 3.1 Global explanation

The simplest explainer uses a single linear factor model with one set of factors to explain for all instances.

$$\tilde{y}_g = \sum_r w_g^{(r)} x^{(r)} \tag{1}$$

where $x^{(r)}$ is the $r$th feature value of the instance with $x^{(0)} = 1$, $w_g^{(r)}$ is the explanation factor for that feature with $w_g^{(0)}$ as the bias term, and $\tilde{y}_g$ is the estimated AI prediction. The Global explanation is trained by fitting a linear regression model on the whole training dataset with mean squared error (MSE) as the training loss against

the AI model's prediction not the ground truth. Fig. 2 shows our user interface (UI) implementation of a Global explanation of how an AI System predicts the price of a house based on 4 attributes, and the bias term which we name as "adjustment".

## 3.2 Subglobal explanation

While the Global explanation is simple for users to understand, its small number of factors limits its expressiveness, so it may not be very faithful to the AI System predictions, i.e., $\tilde{y}_g$ is not close to $\hat{y}$. Instead of adding more complexities for users to interpret, the explanation faithfulness can be increased by partitioning instances into multiple subspaces. Each subspace is then modeled with a separate sparse linear factor explanation. We constrain explanations of each subspace to have the same attributes, and enforce the partition based on binary univariate rules, i.e., inequality on one attribute (e.g., $x_2 \geq 2.5$). Thus, a Subglobal explanation has the form

$$\tilde{y}_s = \sum_\varsigma \sum_r [x \in s_\varsigma] w_{s\varsigma}^{(r)} x^{(r)} = \begin{cases} \sum_r w_{s0}^{(r)} x^{(r)}, & \text{if } x \in s_0 \\ \sum_r w_{s1}^{(r)} x^{(r)}, & \text{if } x \in s_1 \\ \dots \end{cases} \tag{2}$$

where $w_{s\varsigma}$ is the weights of the $\varsigma$th subspace explanation model, $[\cdot]$ is the Iverson bracket that is 1 if its expression is true or 0 otherwise, and $S_\varsigma$ is the set of instances in subspace $\varsigma$. Eq. 2 shows that each subspace has different weights (factors) which it applies to instances within its boundaries.

Training the Subglobal explanation model requires learning the partition boundaries of the subspaces, and the weights of each subspace model. We achieve this by training a linear model tree [13] on the whole training dataset with MSE for the training loss. Such trees are different from common classification decision trees that predict a probability distribution of categorical labels $p(\hat{y})$ at leaves, or regression decision trees that predict a scalar number $\hat{y}$ at the leaves. Instead, linear model trees predict a linear regression model $w_\varsigma$ at each leaf, where each leaf represents a subspace $S_\varsigma$. During training, for each branch in the decision tree, the training algorithm iterates through all features and possible splits, training a linear model for each subspace ($s_\leq$ and $s_\geq$), measuring the combined loss for both models, and choosing the split with the lowest combined
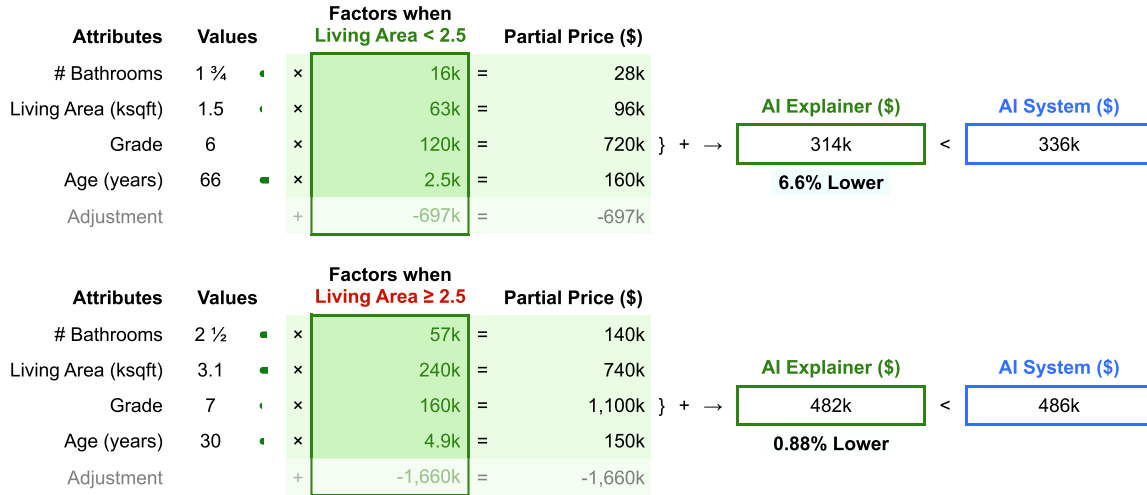
**Figure 3: User interface (UI) of Subglobal explanations for a typical instance (top), and an outlier instance (bottom). Factors are different for each subspace but apply in a fixed way to any instance in each subspace. For example, while small houses with Living Area < 2.5 ksqft have each bathroom being worth \$16k, larger houses have much costlier bathrooms at \$57k.**
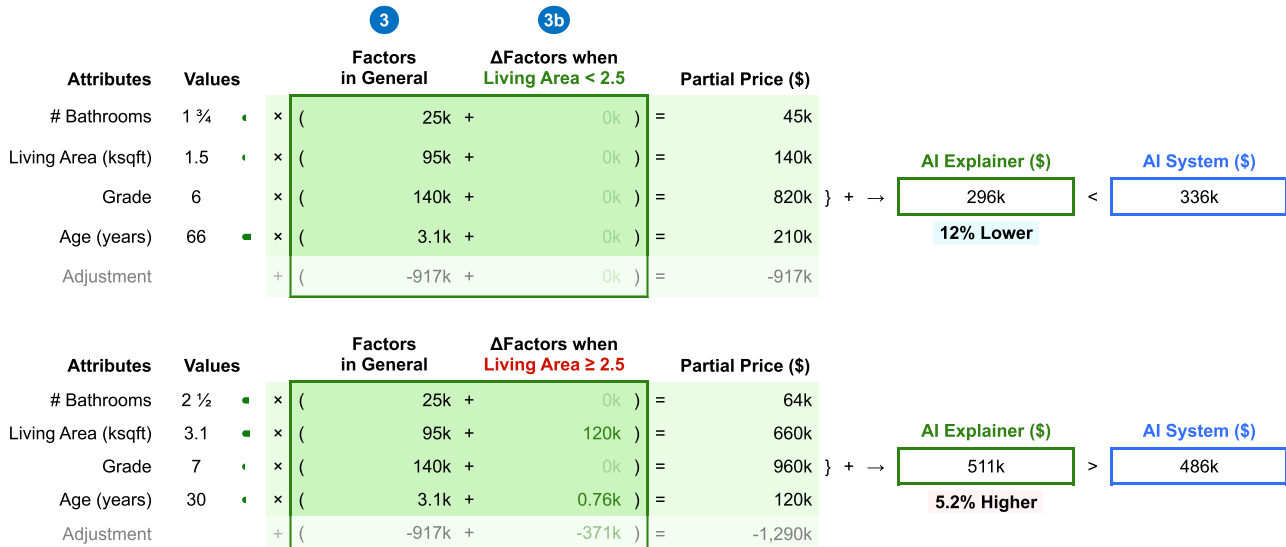


**Figure 4: User interface (UI) of Incremental explanation for an instance in the typical subspace with Living Area < 2.5 (top), and an outlier instance in the minority subspace with Living Area ≥ 2.5 ksqft (bottom). Factors are different for each subspace to fit them accurately. Unlike Subglobal explanations, an additional column (3b) is used to show how factors are incrementally different for the outlier cases. The main factors (3) are the same for both subspaces. For example, while smaller houses have a modest rate of price increase per living area (\$95k/ksqft), larger houses have a rate that is \$120k/ksqft higher (\$215k/ksqft).**

loss. We then assign the majority subspace with the larger dataset as "typical" and minority one as "outliers" (although this can be flexibly adapted to fix user preferences or standard conventions). Though linear model trees are not a novel technique, they are seldom used in explainable AI, and we extend it for Incremental explanations for our technical contribution, described in the next subsection.

Fig. 3 shows our UI of Subglobal explanations with two subspaces: typical (Living Area < 2.5 ksqft) and minority outlier (otherwise). For simplicity, we specifically train a decision stump (one

branch). Each subspace is defined with simple univariate decision boundaries that are easy to interpret. Note that training a logistic regression or linear support vector machines (SVM) would lead to less interpretable decision boundaries, e.g., "16(# Bathrooms) + 120(Grade) < 697", while training on a decision tree would have a more interpretable rule, e.g., "# Bathrooms ≥ 5 and Grade < 5".
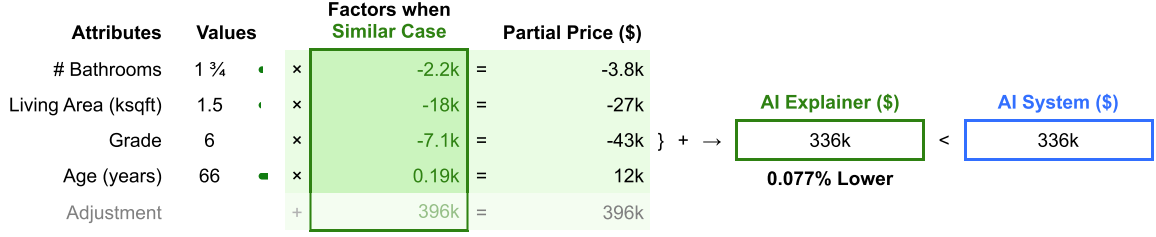
**Figure 5: User interface (UI) of the Local explanation of an instance. Factors are specific to this and similar instances, and will be different for other instances. For example, for houses similar to the one shown, each increase in Grade decreases the house price by $7.1k, but this may not be the case for other houses that have very different attributes.**

## 3.3 Incremental explanation

While Subglobal is more faithful to the AI System than Global, this comes at a cost of the user having to read and remember more factors. The factors are not necessarily consistent between subspaces too, so users would have to interpret them independently. To improve memorability, we propose Incremental explanations that provide general factors for the majority, typical subspace and an incremented factors for special, outlier subspaces. We describe our approach for two subspaces, but it can be extended for multiple subspaces. We define an Incremental explanation as:

$$\tilde{y}_i = \sum_r \left( w_{i0}^{(r)} + \sum_{\varsigma > 0} [x \in s_\varsigma] \Delta w_{\Delta i \varsigma}^{(r)} \right) x^{(r)}$$

$$= \begin{cases} \sum_r w_{i0}^{(r)} x^{(r)}, & \text{if } x \in s_0 \\ \sum_r \left( w_{i0}^{(r)} + \Delta w_{\Delta i1}^{(r)} \right) x^{(r)}, & \text{if } x \in s_1 \\ \dots \end{cases} \quad (3)$$

where $w_{i0}$ is the base factors of the general explanation model, $\Delta w_{\Delta i \varsigma}$ is the incremental factors of the $\varsigma$th special subspace explanation model, $[\cdot]$ is the Iverson bracket that is 1 if its expression is true or 0 otherwise, and $S_\varsigma$ is the set of instances in subspace $\varsigma$. Eq. 3 shows that while the typical subspace has Base factors, all other subspaces have factors defined as an additive Incremental adjustment on the Base factors.

We train the Incremental explanation model in a similar manner as for Subglobal explanations but constrain a dependency in factors across subspaces, i.e., $w_{i\varsigma}^{(r)} = w_{i0}^{(r)} + \Delta w_{\Delta i \varsigma}^{(r)}$. Furthermore, to reduce the number of terms to remember, we aim to keep most incremental weights $\Delta w_{\Delta i \varsigma}$ to be zero. This is achieved by adding a sparsity L1 regularization to the original MSE training loss, i.e.,

$$L(\tilde{y}_i, \hat{y}) = \sum_k \left( \tilde{y}_i^{(k)} - \hat{y}^{(k)} \right)^2 + \lambda \sum_\varsigma \sum_r |\Delta w_{\Delta i \varsigma}| \quad (4)$$

where $k$ indexes the training instances, and $\lambda$ is the regularization hyperparameter. This sparsity regularization makes incremental factors easier to remember, but this trades-off accuracy, so we hypothesize that the Incremental explanation $\tilde{y}_i$ is less faithful than the Subglobal explanation $\tilde{y}_s$. The training algorithm is similar as in Subglobal explanations, but with non-independent parameters for the linear models, extended loss function, and unified optimization of both set of factors. For each candidate split, we set the majority

subspace as "typical" assigning the base factors[1] $w_{i0}$ to it, and specifying Incremental factors $\Delta w_{i\varsigma}$ in other minority "outlier" subspace. Training is performed using gradient descent.

Since the partitioning of subspaces is similar in both Incremental and Subglobal explanations, we keep them the same in the modeling and user studies to avoid partitioning being a confounder. Fig. 4 shows our UI implementation of an Incremental explanation for cases in the typical and outlier subspaces.

## 3.4 Local explanation

We have described sparse linear factors to explain multiple instances, but they can also be used to explain individual instances. Local explanations, such as LIME [51] and SHAP [42], are popular XAI techniques to explain AI decisions on a target instance by training a linear regression model from a dataset of instances that are local (similar) to the target instance. Explaining other instances require retraining other explanation models locally around those instances, and the explanations are not necessarily similar to one another. Hence, Local explanations are faithful to the AI prediction only to instances that are similar, and not globally or subglobally. Consequently, users would need to view many Local explanations to have an overview of the AI behavior across all instances. Although local explanations are not designed for general understanding, their ubiquity encourages their misuse for this objective. We hypothesize that this makes it very difficult for users to estimate how the AI System would predict for new instances or estimate general factors from the inconsistent factors of each instance. We define local linear factor explanations around a target instance $x_l$ as:

$$\tilde{y}_l = \sum_r w_l^{(r)} x^{(r)}, \forall x \approx x_l \quad (5)$$

where $x$ is the instance being explained that is similar to $x_l$, $w_l$ is the weights of the model (factors) with $w_l^{(0)}$ as the bias term, and $\tilde{y}_l$ is the estimation of the local model. We implemented the Local explanation with LIME [51]. Fig. 5 shows our UI implementation of a Local explanation around an instance.

---

[1]Note that the base factors only represent those of the "typical" instances (majority subspace), not of the Global explanation model; this is not the same due to the constraint to minimize incremental factors which will shift the base weights during training.

# 4 EVALUATION

We evaluated Incremental explanations against baseline explanations (Global, Subglobal, Local) across multiple studies to investigate: i) faithfulness to estimate the AI prediction in a modeling study, ii) usage strategies and outcomes to interpret AI decisions in a qualitative formative user study, and iii) impact on decision duration, explanation recall, and AI decision understanding in a quantitative summative user study.

## 4.1 Modeling Study

We conducted a modeling study to evaluate how faithfully each explanation model estimates the AI. We evaluate on three datasets, and our approach can further generalize since we are using standard machine learning processes. We describe the dataset preparation, methods to train and test the models, and evaluation results.

*4.1.1 Applications and datasets.* We evaluated the sparse linear factor explanations on a regression prediction task, since the predictions remain linear, unlike classification that would have tapered effects at high or low probabilities (e.g., logistic regression). Like Poursabzi-Sangdeh et al. [48], we evaluated on a housing price dataset due to the simplicity of the application scenario that most users can readily understand and appreciate. However, we chose not to reuse their NYC dataset since, surprisingly, a linear global model is sufficient to predict prices highly accurately. However, real-world datasets tend to be more complex, and require nonlinear models. Hence, we used the "House Sales in King County, USA" dataset [21] with 21,613 instances to predict the price of houses with 22 features. Prices ranged from $72k to $7.7M (Median = $452k). We performed feature selection to obtain four features (# Bathrooms, Living area, Grade, Age) to limit the cognitive load for users.

For generality, we further evaluate on two additional datasets: Heart Disease [23] with 1025 instances to predict heart disease using 14 common features, and Auto MPG [50] with 398 instances to predict the miles per gallon fuel efficiency using 7 features. Although the prediction task for heart disease is to classify whether a patient has heart disease, the predictor model produces a numeric confidence that can be interpreted as a continuous risk score. We train subsequent explainer models as a regression task to predict the risk score of the predictor model. For the heart disease dataset, to support human interpretability of the explainer models, we performed feature selection to obtain four features (Age, Resting blood pressure, Cholesterol, Max heart rate). Similarly, for the Auto MPG dataset, we performed feature selection to obtain four features (Cylinders, Displacement, Horsepower, Weight).

For simplicity, we partitioned each dataset into two subspaces and set the same rule boundary for both Subglobal and Incremental explanations. The optimal partitions were at Living Area ≥ 2.5 ksqft for House sales, Age ≥ 58 years for Heart disease; and Horsepower ≥ 92W for Auto MPG.

*4.1.2 Results on performance of AI prediction models.* For each dataset, we trained a random forest regressor (House Sales, Auto MPG) or classifier (Heart Disease) [9] as an AI prediction model on a training set of 80% instances, and evaluated on a heldout test set of 20% instances. We then trained the four explanation (XAI) types

to explain all instances in the test set. For Local explanation, we averaged the performance across all instances. Using the training set, we performed 5-fold cross validation in all our analyses and report the mean and standard deviation of the validation performance averaged across folds; see Table 2. We report the performance on the heldout test dataset – House Sales: mean absolute error (MAE) = $139k and $R^2$ = 0.67; Heart Disease: accuracy = 86% and test AUC = 0.86; Auto MPG: MAE = 3.12 mpg and $R^2$ = 0.71.

*4.1.3 Results on faithfulness of XAI types.* Fig. 6 shows the unfaithfulness of the XAI types calculated by their absolute error (AE) between the explainer $\tilde{y}$ and predictor predictions $\hat{y}$. Note that the explanation faithfulness ($\tilde{y}$ vs. $\hat{y}$) does not measure the same thing as predictor performance ($\hat{y}$ vs. $y$). As expected, Global explanations had the worst faithfulness due to its low expressiveness (fewest factors), while Local explanations were the best because they were trained to be accurate to small local neighborhoods. However, they are not robust or memorable, so we expect users to not gain as much understanding from them compared to the other XAI types. Subglobal explanations had slightly better faithfulness than Incremental explanations since the latter had another objective of simplicity (fewer incremental factors). However, we expect this difference to be negligible in practical use by people, and the problem of memorability or cognitive load would override the small benefit of faithfulness, and we hypothesize that Subglobal explanations are less memorable and interpretable than Incremental explanations. We evaluate these hypotheses later in the summative user study.

*4.1.4 Results on performance of XAI as glassbox explainers.* We further evaluate whether Subglobal and Incremental explanations can serve as accurate interpretable "glassbox" models. In such cases, these models would be used for the AI prediction task, and be intrinsically interpretable, thus avoiding any unfaithfulness of surrogate explanation models. We trained the models on the training dataset and report their performance. Since Heart Disease is a classification task, we accordingly apply logistic activation to the linear regression outputs and change the training objective to the binary cross-entropy loss. The interpretability of the factor coefficients is affected, since the sigmoid transform in logistic regression applies a nonlinear distortion on all weight. However, the directionality and the magnitude of the factors still provide more interpretability than blackbox models. We report the glassbox explainer performance as "AI Performance" in Table 2. In summary, Subglobal and Incremental models performed better than Global models, but this is still a worse than the nonlinear AI model (random forest, in this case).

*4.1.5 Investigating explanations for multivariate attributes across subspaces.* In Fig. A.1 in Appendix A.1, we show the 2D decision surfaces of the AI and explanation models for the House Sales dataset to demonstrate how the four different XAI types model the relationships between the AI prediction and multivariate attributes.

*4.1.6 Investigating varying subspace thresholds.* While Subglobal and Incremental explanations learn the feature space partitioning threshold automatically with the linear model tree, they can also be set with custom values to fit the explanation needs. We thus examine, in Appendix A.2, how selecting different partition thresholds affect the faithfulness of each subspace explainer model, how the factors change, and whether incremental factors are kept small.

**Table 2: Modeling results from 5-fold cross-validation of AI performance and XAI faithfulness across three datasets showing mean ± standard deviation. AI performance indicates when an explainer is trained on the *ground truth* dataset as a glassbox interpretable model. XAI unfaithfulness evaluates each explainer as a surrogate explanation with respect to the AI Model. Except for AI performance for Heart Disease that is measured as % Accuracy, all other metrics are MAE, where smaller is better.**

| House Sales | | AI | XAI types | | | |
|---|---|---|---|---|---|---|
| Subspace | Metric (MAE $k) | Model | Global | Subglobal | Incremental | Local |
| Combined | AI Performance (inv) | 132.5 ± 2.4 | 145.1 ± 3.0 | 138.1 ± 2.6 | 139.8 ± 2.9 | – |
| | XAI Unfaithfulness | 0 | 68.4 ± 1.2 | 48.5 ± 0.7 | 53.8 ± 0.7 | 32.4 ± 0.2 |
| Typical | AI Performance (inv) | 102.9 ± 1.3 | 113.5 ± 1.5 | 105.5 ± 1.5 | 108.5 ± 1.6 | – |
| | XAI Unfaithfulness | 0 | 54.3 ± 1.5 | 35.1 ± 0.5 | 41.3 ± 1.1 | 26.9 ± 0.4 |
| Outlier | AI Performance (inv) | 213.2 ± 7.9 | 231.2 ± 8.9 | 227.0 ± 9.2 | 225.1 ± 9.0 | – |
| | XAI Unfaithfulness | 0 | 107.0 ± 3.4 | 85.1 ± 2.7 | 87.7 ± 2.1 | 47.0 ± 1.8 |

| Heart Disease | | AI | XAI types | | | |
|---|---|---|---|---|---|---|
| Subspace | Metric (Acc %, MAE %) | Model | Global | Subglobal | Incremental | Local |
| Combined | AI Performance | 85.4 ± 2.7% | 69.87 ± 6.66% | 71.46 ± 3.06% | 70.6 ± 4.55% | – |
| | XAI Unfaithfulness | 0 | 18.2 ± 0.8 | 15.5 ± 0.9 | 15.7 ± 1.1 | 8.9 ± 0.5 |
| Typical | AI Performance | 85.7 ± 2.2% | 79.73 ± 3.71% | 78.3 ± 3.5% | 77.88 ± 2.57% | – |
| | XAI Unfaithfulness | 0 | 16.9 ± 1.1 | 15.3 ± 1.5 | 15.3 ± 1.6 | 8.5 ± 0.9 |
| Outlier | AI Performance | 84.9 ± 4.7% | 55.85 ± 11.05% | 61.55 ± 4.68% | 60.2 ± 7.56% | – |
| | XAI Unfaithfulness | 0 | 20 ± 1.7 | 15.8 ± 0.6 | 16.5 ± 0.6 | 9.4 ± 0.9 |

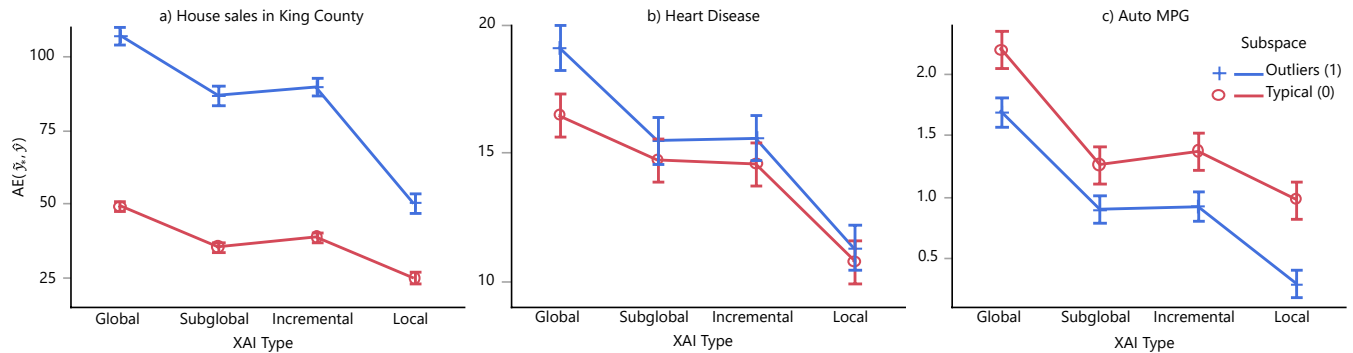| Auto MPG | | AI | XAI types | | | |
|---|---|---|---|---|---|---|
| Subspace | Metric (MAE mpg) | Model | Global | Subglobal | Incremental | Local |
| Combined | AI Performance (inv) | 2.65 ± 0.3 | 3.25 ± 0.32 | 2.78 ± 0.28 | 2.84 ± 0.34 | – |
| | XAI Unfaithfulness | 0 | 1.91 ± 0.16 | 1.08 ± 0.05 | 1.11 ± 0.09 | 0.53 ± 0.04 |
| Typical | AI Performance (inv) | 3.37 ± 0.42 | 3.92 ± 0.49 | 3.59 ± 0.51 | 3.64 ± 0.53 | – |
| | XAI Unfaithfulness | 0 | 2.17 ± 0.25 | 1.32 ± 0.14 | 1.36 ± 0.14 | 0.83 ± 0.13 |
| Outlier | AI Performance (inv) | 1.93 ± 0.58 | 2.58 ± 0.30 | 1.97 ± 0.49 | 2.03 ± 0.46 | – |
| | XAI Unfaithfulness | 0 | 1.69 ± 0.30 | 0.87 ± 0.17 | 0.90 ± 0.21 | 0.23 ± 0.05 |



**Figure 6: Results of modeling study showing the unfaithfulness of each explanation type calculated as absolute error (AE) between the AI Explainer estimation and AI System prediction $AE(\tilde{y}, \hat{y})$ across three prediction tasks with different datasets: a) House Sales in King County [21], b) Heart Disease [23], c) Auto MPG [50]. Global explanations are least faithful, Local explanations most faithful, and Subglobal and Incremental explanations have similar faithfulness. The faithfulness of typical or outlier cases depends on the explainer models trained for each dataset.**

## 4.2 Formative User Study

To investigate how people use the XAI types (Global, Subglobal, Incremental, Local) and identify usability issues, we conducted a formative study with 14 participants recruited from a local university. They were 23 years old on average (20 to 30), and 6 were female. All were undergraduate and graduate students from various disciplines (5 Sciences, 4 Business, 2 Engineering and Technology, 2 Arts, and 1 Healthcare). The study was conducted virtually over a Zoom video call with screen recording, and lasted 60 minute. Participants were compensated with a digital payment of $15 USD.

*4.2.1 Method and Procedure.* We conducted the study with a *within-subjects* experiment design, where each participant views multiple XAI types, so that they may directly compare among them. The experiment apparatus and procedure are similar to the subsequent summative study, which we describe later. To ensure that participants did not confuse between attribute names, values, factors, and

partial contributions, we trained them to distinguish the columns in the tabular UI, understand how each partial contribution is calculated as a multiplication of factor (weight) and value, and verified their understanding with screening questions. Each participant first used the Global explanation as a baseline, then 1-3 randomly selected XAI types as time permitted. For each explanation, the participant performed 3 trials of viewing an AI explanation to predict the price of a house instance. They were asked to estimate what they thought the AI System would predict based on the information provided by the explanation. The instances were chosen from the same dataset as the Modeling study, and used the same apparatus as the summative study (conducted later), with user interfaces shown in the Technical Approach section.

We used the think aloud protocol to elicit the participant's thought processes as they read and applied the explanations. The participant could ask clarification questions any time too. Since the participants were guided and supported by the experimenters, we do not report the performance of their estimations on the AI system. With participant consent, we audio and screen recorded participant vocalizations and interactions with the UI.

*4.2.2 Findings.* We conducted a thematic analysis on participant behaviors and report key findings.

We note that some participants may have conflated the AI and XAI behavior, but we do not require our users to treat them as separate, since both are meant to be presented as a unified agent in the AI's user interface. In this study, we focus on how participants interpreted each XAI type, rather than their trust or decision with respect to the AI prediction model.

*a) Dynamic explanations perceived as more realistic than static, global explanations.* Most participants preferred explanations to be dynamic rather than static like in Global explanations. Only P13 felt that *"the fixed factors of [Global] are more intuitive and similar to how many people think."* Perhaps, she preferred rules of behaviors to be consistent and unchanging. On the other hand, many participants appreciated the complexity of house price estimations and AI systems. P1 believed that the AI system *"is a bit more dynamic in nature, or the equations will adjust accordingly to how much data is set into the thing."* He felt that Global explanations did not reflect the AI system well since *"it would just come up with one static figure because the factors itself is consistent and doesn't change across the house type."* P7 expected to see different factors for instances of different categories, remarking that *"it's not very realistic for the factors to be the same for different house types, like factor for bathrooms is always [the same]"*. This suggests she categorizes instances into types and expect rules to apply different for each category. In contrast, P8 appreciated the adaptiveness of non-Global explanations and felt that *"it's logical that the factors would change for different type of houses, ... since there might be other factors that influence the factor values for each attribute."* Similarly, on seeing that *"all the factors were the same"* for the Global explanation, P6 remarked *"that might not be good."* He explained that for larger houses, the factor for Living Area *"should be on a diminishing graph"*, i.e., smaller factor than for smaller houses due to smaller marginal utility of living area in an already large house.

Though less dynamic than Local explanations, participants found the subspace partitioning of Subglobal and Incremental explanations intuitive. P1 explains, *"[Subglobal] is a lot more accurate [than Global] because it considers more things"*, referring to the two sets of factors given. P14 affirmed that *"the additional factors [in Incremental] are helpful for the predictions in terms of accuracy"*. P3 remarked that *"[the Incremental factors] makes sense, because for bigger houses the land would cost more."* This also shows the relative understanding that P3 had to compare between subspaces, thus demonstrating the usefulness of explaining incrementally.

*b) Incremental explanations perceived as more memorable and accurate than Subglobal explanations.* Both Incremental and Subglobal explanations partition the subspaces similarly, but Incremental articulates the relationship between the two subspaces, and Subglobal treats them independently. Participants could appreciate the benefit of providing this context in Incremental explanations. P6 liked that *"[Incremental] would be more informed since you are telling the user how they're changing the factors, that there is an addition"*. P11 even stated that the consistency of Incremental made him feel assured because *"not all the factors are changing, like there were more considerations being made by the explanation."* P4 mentioned that *"[Incremental] would be easier for me to remember because there are fewer numbers"* and P9 agreed that this is due to *"rather than remembering two separate sets of factors"*. Furthermore, P12 believed *"[Incremental] will give you more accurate values, which helps you make decisions quicker."* Though, this is not necessarily true, and suggests a positive halo effect of better usability leading to perceived correctness. Nevertheless, it suggests that this can help boost user confidence, trust and usage of Incremental explanations. Despite these benefits, some participants faced some usability issues. P7 felt that *"[Incremental] feels logical... but more time-consuming since it's slightly complex due to the additional factors you have to add for the calculation."*

## 4.3 Summative User Study

We conducted a summative user study to evaluate the interpretability and memorability of each XAI type. We investigate how well participants understand, remember, and apply explanations to anticipate behavior for future instances. While testing the impact on a downstream decision making task would be meaningful, it would impose experiment confounders, such as the participant's prior knowledge of the task [38], their varying underlying utility objectives (e.g., how much they care about cheap housing) [43], increased mental fatigue which limits the number of trials [2], and conflation between AI and XAI estimations. Thus, we leave that for future work.

Next, we describe our experiment design and hypotheses, experiment apparatus, procedure, analysis and results.

*4.3.1 Experiment Design.* We designed our experiment as a 4×2 factorial mixed-design experiment with primary independent variable (IV) as **XAI type** (four levels: Global, Subglobal, Incremental, Local) and secondary IV as **Subspace** segment (two levels: typical, special) to investigate if effects differ by instance type. XAI type was manipulated between-subjects due to the learning effect of participants sticking to one mental model of the AI Explainer after being trained on the first XAI type. Subspace was manipulated

**Table 3: Hypotheses and findings of the summative user study regarding different dependent variables for various XAI types: Global (G), Subglobal (S), Incremental (I), Local (L).**

| Dependent Variable | Metric | Hypothesis | Finding | Evidence |
|---|---|---|---|---|
| Decision duration | Log(Time) | G < I < S < L | I ≈ G < S ≈ L | Fig. 9a |
| Explanation evocation | $-AE(\tilde{y}_h , \tilde{y})$ | L < S < I < G | I ≈ S < L ≈ G | Fig. 9b |
| Supported understanding | $-AE(\hat{y}_h, \hat{y})$ | G < I < S < L | G < S ≈ L ≈ I | Fig. 9c |
| Sustained understanding | $-AE(\hat{y}_h \mid \tilde{y}, \hat{y})$ | L < G ≈ S ≈ I | G < L ≈ I ≈ S<br>(Special: G < L < S ≈ I) | Fig. 9d |
| Explanation recall | $-AE(w_h^{(0)}, w^{(0)})$ | L < S < I < G | L ≈ G ≤ S ≈ I<br>(L ≈ G ≤ I) | Fig. 10 |
| Perceived helpfulness<br>Perceived ease-of-task | 7-pt Likert scale | L < G < S < I | L ≈ G ≈ S ≈ I | Fig. 11 |

within-subjects by selecting 100 instances from the full datasets, where we balanced 50 typical and 50 special. Each participant is tested on 30 randomly selected instances.

We measured several objective dependent variables to evaluate explanation recall, application, and understanding:

- *Explanation recall* measures how accurately the participant can infer or remember each factor $w^{(r)}$ of the XAI type, by typing them out. This explicitly measures *memorability*. For the $r$th attribute, given the participant's estimate $w_h^{(r)}$, we calculate the lack of recall by the MAE. We asked about the factors for all instances (global), typical or special instances (subglobal). Although participants with Local explanations never see general explanations, we ask them to infer broadly.
- *Sustained understanding (without XAI)* measures how well the participant can estimate the AI System's prediction. This is *forward simulatability* [16], a popular metric in XAI research and evaluations. Since we are modeling a regression problem (rather than the typical classification), we calculate this with a proxy metric for unfaithfulness with the absolute error (AE) of regression predictions, i.e., $|\hat{y}_h - \hat{y}|$. This measures how well the participant can apply knowledge gained from studying explanations for other instances without having seen their explanations. It measures deeper understanding than Supported understanding, which we also measure, described next.
- *Supported understanding (with XAI)* measures how well the participant can estimate the AI System's prediction, given that he/she can view an approximation from the AI Explainer $\tilde{y}$, i.e., $\hat{y}_h|\tilde{y}$. This is similar to Sustained understanding, but easier, since the participant can leverage $\tilde{y}$ to estimate his answer.
- *Explanation evocation* measures participant correctness to estimate the AI Explainer's estimation $\tilde{y}$. This is the forward simulatability of the AI Explainer, which we compute its reverse as $\tilde{y}_h - \tilde{y}$. Unlike Explanation recall, which directly elicits explanatory factors, this queries the participant about the explanation outcome, which implicitly evokes the explanation.
- *Decision duration* measures how long participants spent to perform the forward simulatability task without XAI. Since duration follows a long-tail distribution, we analyzed its logarithm.

For participant convenience, we measured numeric factors and prediction estimates with sliders to give users bounds when answering, but use the same wide range to avoid priming. Furthermore, we measured subjective opinions on *Perceived helpfulness* and *Perceived ease-of-task* to investigate how helpful the different XAI types were. Specifically, we asked whether the participant agreed or disagreed that: the AI was accurate (1), the explanation was helpful to estimate factors globally (2) and subglobally (3), the forward simulatability tasks were easy with (4) or without (5) the explanation. These were measured on a 7-point Likert scale (-3 to +3). Table 3 summarize our hypotheses and the subsequent findings from our results and analysis, described later.

*4.3.2 Experiment Apparatus.* Our user interface was inspired by the linear factors explanation interface of Poursabzi-Sangdeh et al. [48], but we adapted it to distinguish the linear model explainer from the nonlinear model predictor, and extended it to support various XAI types: Global, Subglobal, Incremental, Local. Participants saw the exact interface as shown earlier in Figs. 2-5. See the Appendix for the full survey that participants saw. We also made the UI interactive (see Fig. 7) to facilitate participant learning and engagement by examining how explanations and predictions depend on instance values and factors. To improve interpretability and usability, we rounded most numbers to two significant figures, though the calculations are still done in full precision, and participants can see the precise numbers by hovering their mouse cursor. The intercept term is rounded to three significant figures, since it is a direct partial contribution component, unlike the other terms as factors. We included green meter bars to show the relative levels of attribute values to allow participants to interpret the sense of each number. During training, we show the % error between the AI Explainer and AI System to make this salient and accelerate learning about explanation faithfulness. We implemented our survey in Qualtrics and embedded the user interface.

*4.3.3 Experiment Procedure.* Each participant was engaged in the following procedure:

1) Introduction to the study (see Appendix Figs. A.3-A.5).

**Figure 7: User interface (UI) during testing with factors hidden, but editable. Participants can type their own numbers to explore how the AI Explainer would compute based on various factors. This helps users to learn how factors work. Here, participants are asked to forward simulate both the AI Explainer and AI System without seeing any factor explanations.**

2) Consent to participate. This study was approved by the university institutional review board (IRB).

3) Tutorial on the AI prediction task (housing price prediction) (see Fig. A.6-A.8).

4) Tutorial on the user interface and test tasks. Different features are introduced depending on XAI type condition (see Figs. A.10, A.12, and A.14).

5) Screening questions to ensure that the participant can interpret and use the explanation factors correctly (see Figs. A.9, A.11, A.13, and A.15).

6) Forward simulatability session (see Fig. 8) of 5 trials with reflection and 25 regular trials, where each trial:

i) On page 1 (see Fig. A.16), *view* the user interface with only values shown and *forward simulate* what the AI Explainer and AI System will output, as *Explanation evocation* ($\tilde{y}_h$) and *Sustained understanding* ($\hat{y}_h$), respectively. Factors and consequent calculations are hidden. The UI is interactive to allow the participant to type different factor values while attempting to estimate the AI outputs (see Fig. 7). Participants used two sliders to indicate their estimates for the AI Explainer and AI System outputs (see the bottom of Fig. A.16). To enhance the learning of the AI Explainer and AI System behavior and help participants learn to apply the factors, we posed several reflection questions (see the middle section of Fig. A.16). These were only asked for the first 5 trials to limit the survey duration.

ii) On page 2 (see Fig. A.17), additionally *view* explanatory factors and AI Explainer calculations, based on XAI type condition, and *forward simulate* what the AI System will predict. This measures *Supported understanding*.

iii) On page 3 (see Fig. A.18), *review* their answers with the actual AI System prediction. This frequent review allows the participant to continuously learn from his/her mistakes to pay better attention to learn the factors, and strengthen their understanding, across all conditions.

7) Factors recall session (see Figs. A.19-A.21), to recall the factors for a) all instances in general, b) typical instances, and c) outlier instances. No specific instance values are shown, but the participant can enter their own values to examine what factors could be suitable.

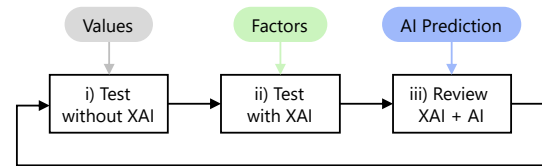8) Answer ratings questions on *Perceived helpfulness* and *Perceived ease-of-task*.



**Figure 8: Procedure of a trial in the summative user study to evaluate explanation understanding and memorability.**

9) Answer demographics questions.

10) Acknowledge bonus calculations and exit.

We provided an incentive bonus of £0.03 for each Sustained understanding task if the participant could estimate the AI System prediction correctly to within 10% relative error (max £2.70), and max of £0.15 for each Explanation recall task (3 tasks) based on the mean relative error $\varepsilon_{MRE}$ on a test set of 100 instances, calculated as £0.15 × $(1 - \varepsilon_{MRE})$.

*4.3.4 Participants.* We recruited workers from Prolific.co, where 160 passed screening and 336 failed. Participants who completed the study had an median age 37 years old (26 to 81), and were 35% female. Participants completed the survey in a median time of 96 min, and were compensated with a base of £9.00 and Median bonus of £1.20 (£0 to £2.88).

*4.3.5 Statistical Analysis.* We performed a linear mixed effects model fit on each dependent variable as the response, XAI type and Subspace, along with other confounding variables as fixed effects, some interaction effects among the factors, and Participant as random effect. See Table 4 for details. Note that Supported understanding and Sustained understanding are calculated from the same measure, forward simulatability, but differ only by when the task was posed, before and after showing XAI, respectively. Since they share the dependent variable, we analyze these responses with a single linear mixed effects model with Test with XAI to distinguish between the two types of understanding.

The model fit was good for Log(Task time) ($R^2$ = .628) indicating that task time depended much on fixed effects XAI type, subspace, trial sequence specific test instance, and the participant random effect. The model fit for Explanation recall was good ($R^2$ = .763), indicating that it was influenced much by two factors (XAI type

**Table 4: Statistical analysis of responses due to effects (one per row), as linear mixed effects models with random effects, fixed effects, and their interaction effect. $F$ and $p$ values indicate ANOVA tests and $R^2$ indicate model goodness-of-fit.**

| Response | Linear Effects Model (Participants as random effects) | F | p>F | $R^2$ |
|---|---|---|---|---|
| | XAI Type + | 13.3 | <.0001 | .483 |
| | Subspace + | 139.8 | <.0001 | |
| Explanation evocation | XAI Type × Subspace | 20.0 | <.0001 | |
| | Trial ID + | 127.8 | <.0001 | |
| | Instance ID | 7.5 | <.0001 | |
| | XAI Type + | 18.2 | <.0001 | .386 |
| | Subspace + | 242.0 | <.0001 | |
| | Test with XAI + | 912.1 | <.0001 | |
| | XAI Type × Subspace | 12.2 | <.0001 | |
| Supported and Sustained Understanding | Test with XAI × Subspace + | 28.5 | <.0001 | |
| | XAI Type × Test with XAI + | 22.6 | <.0001 | |
| | XAI Type × Test with XAI × Subspace + | 4.5 | .0039 | |
| | Trial ID + | 94.6 | <.0001 | |
| | Instance ID + | 15.2 | <.0001 | |
| | XAI Type + | 5.8 | .0006 | .628 |
| | Subspace + | 24.8 | <.0001 | |
| Log(Task time w/o XAI) | XAI Type × Subspace + | 5.8 | .0006 | |
| | Trial ID + | 4854.0 | <.0001 | |
| | Instance ID | 1.5 | .0016 | |
| | XAI Type + | 8.6 | <.0001 | .763 |
| Explanation recall | Subspace + | 34.5 | <.0001 | |
| | XAI Type × Subspace | 1.7 | n.s. | |
| Perceived helpfulness | XAI Type | 1.7 | n.s. | .977 |
| Perceived ease-of-task | XAI Type | 0.1 | n.s. | .863 |

and Subspace). The model fits for Perceived helpfulness and ease-of-task were also very good ($R^2$ = .977 and .863, respectively), though there were no significant effect due to XAI type, suggesting high variance based on participant individual effect. The model fit was slightly poorer for Explanation evocation ($R^2$ = .483), due to the difficulty to recall the explanation factors and apply weight sum arithmetic to estimate the AI Explainer's prediction $\tilde{y}_h$, leading to increased variance in participant performance. The model fit for Understanding was somewhat low ($R^2$ = .386), because estimating the AI System's prediction regardless of explanation ($\hat{y}_h$ and $\hat{y}_h|\tilde{y}$) are even more difficult and uncertain than estimating $\tilde{y}_h$. When analyzing Supported and Sustained understanding in separate models instead of a larger model with the "Test with XAI" factor, we obtained better model fits ($R^2$ = .415 and .479, respectively) which is similar as for Explanation evocation, but this does not properly account for viewing XAI as a causal factor.

*4.3.6 Quantitative Results.* Table 4 summarizes the model fits in terms of our hypotheses. We describe our results in terms of each dependent variable and summarize the findings with respect to each XAI type. All fixed effects reported are very statistically significant (p<.0001), and describe specific comparisons based on contrast tests. We discuss i) how well participants could estimate the AI Explainer and AI System predictions given different XAI types based on the Forward simulatability trials, ii) their ability to recall the factors of each XAI type in the Factors recall session, and iii) their perceptions on XAI helpfulness and ease-of-task.

*Forward simulatability tasks.* Participant performance varied across XAI types, but were generally poorer for special than typical

cases (p<.0001). See Fig. 9. Next, we discuss specific effects and interpret their effect sizes[2].

a) *Decision duration:* Participants who were trained on Global or Incremental explanations were 1.19 (95% CI: 1.10 to 1.30) times[3] faster ($M_{G,I}$ = 46.8s vs. $M_{S,L}$ = 55.7s) at determining the AI System's output than those trained on Subglobal and Local explanations (Contrast test: ΔLog(Time) = 0.173 ± 0.045, p<.0001).

b) *Explanation evocation:* Participants with Incremental or Subglobal explanations were more accurate in estimating the AI Explainer's output by $97.5k ± $34.6k (95% CI) than those with Global or Local explanations ($M_{S,I}$ = $134.2k vs. $M_{G,L}$ = $231.6k, contrast test p<.0001). Given the average house price of $589.8k, this is 16.5% lower error.

c) *Supported understanding:* Participants who viewed Global explanations were worst (highest AE) at estimating the AI System output by $79.1k ± $13.0k (95% CI), with 13.4% lower error, *even after* viewing the AI Explainer's output than those who viewed Subglobal, Incremental, or Local ($M_G$ = $145.1k vs. $M_{S,I,L}$ = $66.0k, contrast test p<.0001).

d) *Sustained understanding:* The trends in participant performance was similar as for Supported understanding, but were worse due to the increased difficulty of estimating without first viewing

---

[2]Due to how variance is calculated in linear mixed effects models [54], there is no agreed way to calculate standardized effect sizes for fixed main or interaction effects. Hence, we report unstandardized effect sizes, which are the raw differences of the response variables. These are adequate to convey the practical significance in application-specific contexts and interpret their practical meaningfulness [17].
[3]Calculated by inverse transforming ΔLog(Time), i.e., $\exp(\log(t_2) - \log(t_1)) = \exp(\log(t_2/t_1)) = t_2/t_1$.
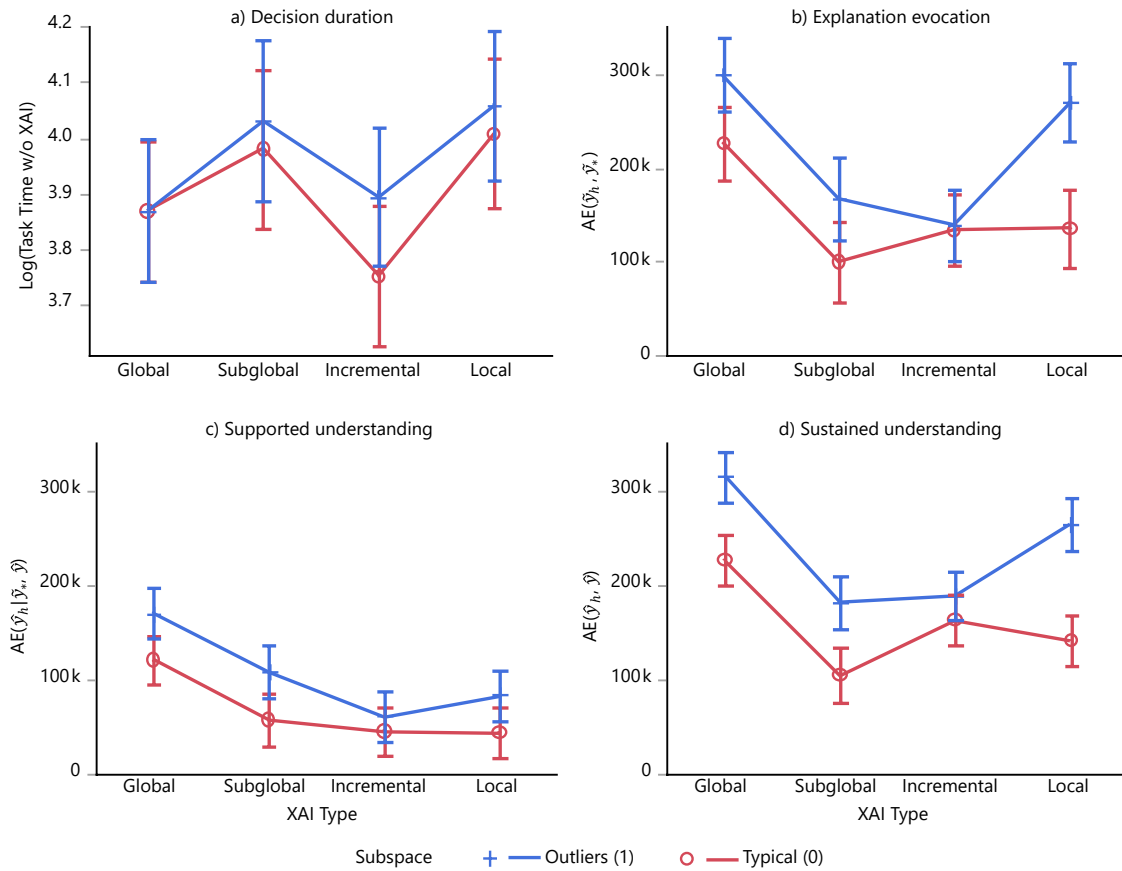
**Figure 9: Results from forward simulatability trials to estimate the AI Explainer and AI System outputs without viewing explanations (b, d), and estimate the AI System output with explanation with timing (a, c). Error bars indicate 90% confidence interval.**

AI explanations. Participants who were trained on Global explanations were worst (highest AE) by $97.3k ± $13.0k (95% CI), with 16.5% lower error, compared to those trained on other explanation types ($M_G$ = $270.1k vs. $M_{S,I,L}$ = $172.6k, contrast test p<.0001). Furthermore, for Special cases, participants trained on Incremental and Subglobal explanations were better by $104.6k ± $13.0k, with 17.7% less error, than those trained on Local explanations ($M_{S,I}$ = $104.6k vs. $M_{G,L}$ = $292.8k, contrast test p<.0001); this suggests that subspace explanations help users to better understand special cases.

*Explanation recall task.* We analyzed how well participants could recall or infer each factor for any instance in general (globally) or for typical or special cases (subglobally). While recall for most factors across global/subglobal were not significantly different, the recall for the explanation intercept term $w^{(0)}$ was notable. Fig. 10 shows the results of recalling $w^{(0)}$ for factor recall sessions of any, typical and outlier cases. Participants recalled factors from Incremental explanations significantly better (lower AE) by $456k ± $185k (95% CI) than from Global and Local explanations (contrast test p<.0001), which is practically significant compared to the intercept terms

−$1,040k (Combined), −$697k (Typical), −$1,660k (Outliers). Furthermore, though recalling Incremental factors was slightly better than of Subglobal factors, this was not significant (p = n.s.).

*Perception ratings.* We had posed multiple questions on Perceived helpfulness and Perceived ease-of-task, but found that all perception questions except ease-of-task without explanation were correlated. Thus, we averaged them into a Perceived helpfulness metric (Cronbach's $\alpha$ = .805). Fig. 11b summarizes the results of the perception measures. Participants perceived all XAI types as somewhat helpful (M=0.84 on a -3 to +3 Likert scale), but found the forward simulatability task without XAI somewhat difficult (M=-0.75). There were no significant differences across XAI types.

*4.3.7 Summary of results.* We now summarize our results of how each XAI type compares to others.

- *Global explanation* was among the fastest type due to its simplicity, but it was not the most objectively helpful for understanding due to its low faithfulness.
- *Local explanation* supports better understanding when provided, but this understanding was not sustained for instances without explanations, since participants were unable to learn to infer factors for new cases.
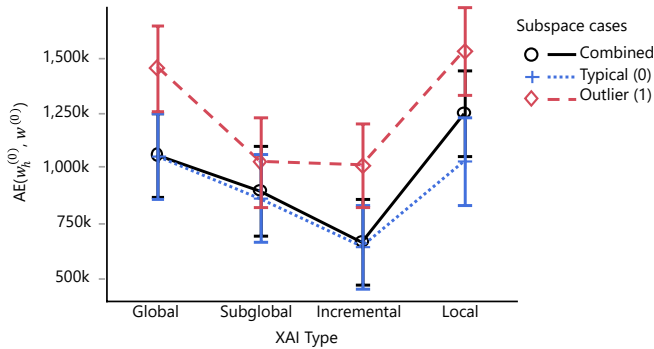
**Figure 10: Results of explanation recall of the intercept term $w^{(0)}$ for the global and subglobal test sessions.**
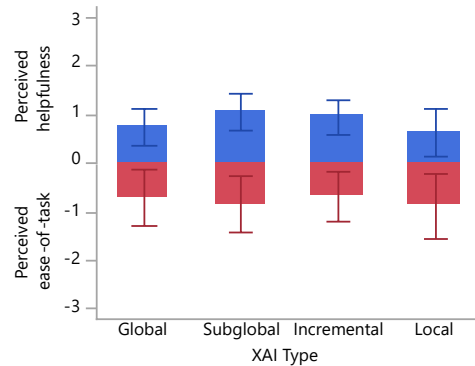


**Figure 11: Results of Perceived helpfulness of AI and XAI, and Perceived ease-of-task without XAI on a 7-point Likert scale from Strongly Disagree (-3) to Strongly Agree (+3). Error bars indicate 90% confidence interval.**

- *Subglobal explanation* supported better recall and understanding than Global or Local explanations, but participants were slow when using them to estimate what the AI System will predict.
- *Incremental explanation* was the fastest to use (as fast as Global explanation), and best for Supported and Sustained understanding (equally good as Subglobal).

## 5 DISCUSSION

We have introduced the paradigm of Incremental XAI, implemented its capabilities and validated its usefulness to help user understanding and recall. Here, we discuss its generalization and limitations.

### 5.1 Generalizing incremental linear factors explanations

Our implementation of Incremental explanations only had a partition along one attribute to divide instances into two subspaces. Nevertheless, since we used a tree-based partitioning method, our approach can apply to more splits and splits in multiple attributes. The splits can also be done on categorical attributes where each subspace can be defined by individual or a set of labels. However, adding more splits will add complexity to the user interface and more information for users to learn and understand, especially in a short online study. Future work is needed to investigate this. Furthermore, we had partitioned our subspaces with trees, but rules may be used to allow fewer terms instead. While tree-based subspaces cannot overlap due to the hierarchical execution of rules, rule sets for different subspaces may overlap, i.e., multiple rules with different features may be true [29]. This can be mitigated with a tie-breaker [29] or by using prioritized rule lists [32].

We had only investigated instances with numeric features, since we focused on training linear factors for explanation models. To accommodate non-numeric features, such as categorical features, standard approaches to convert them, such as one-hot encoding, could be applied. Each categorical level would be interpreted as a feature that is either present (1) or absent (0) with linear factor that is only applied when the level is present.

While our approach reduced the number of factors of the incremented weights, the base weights can also be simplified. This way, users can view an initial explanation with very few attributes, then

incrementally learn new attributes for special cases or further details. This can be accomplished by applying sparsity regularization to the base factors, and a loss penalty to adding new incremental factors (which facilitates new factors if beneficial to the loss). We had strongly limited the number of features to four shown to participants to manage their cognitive load. However, applications in machine learning could to involve about 100 features. Incremental XAI that gradually shows more features can help users to eventually learn many features, and gain an understanding that will be highly faithful to the AI model. This could be implemented by applying modeling on one subspace, iterating >2 steps, and regularizing against reusing features across steps. Further work is needed to model and evaluate on datasets with many features. Though, user testing would be challenging in lab or online studies, since learning new features is harder than adapting prior knowledge about existing features, and learning many features may not be feasible in short durations.

The partitioning of subspaces was determined in a data-driven manner with the tree model, but the factors can still be unwieldy. We had rounded the numeric factors for simplicity, but they could also be constrained as integers or multiples of integers [61]. The split levels and factors could also be relatable [68] and presented verbally or narratively [53], so that users can make sense and better remember them.

### 5.2 Generalizing incremental explanations

We have investigated incremental explanations for linear factor explanations, but argue that this can be generalized to other explanation techniques, such as generalized additive models (GAM) and rules or decision trees. These models can be used for Global explanations by training on the full dataset, or for Subglobal or Incremental explanations by training on subspaces, or Local explanations by training on local neighboring instances. For example, a base function could describe a quadratic trend in a feature, while an incremental function could describe a suppressing cubic effect for special cases; or typical cases could be described with a rule of two features, while special cases could be described by substituting the second feature with a third one for a new rule.

First, we generalize to nonlinear models with independent features, such that each feature $x^{(r)}$ has a partial contribution $f^{(r)}$ to the prediction, i.e., $\tilde{y} = \sum_r f^{(r)}$. We had modeled the contribution of each feature by a linear factor, i.e., $f^{(r)} = w^{(r)}x^{(r)}$ However, but this contribution could be nonlinear, i.e, $f^{(r)} = f^{(r)}(x^{(r)})$. Indeed, this matches the form of GAMs that combine nonlinear effects of features additively, i.e., $\tilde{y} = \sum_r f^{(r)}(x^{(r)})$. Hence, for nonlinear models, extending Eq. 2, a generalized, nonlinear Subglobal explanations is

$$\tilde{y}_s = \sum_{\varsigma} \sum_r [x \in s_{\varsigma}] f^{(r)}_{s\varsigma} \qquad (6)$$

where $f^{(r)}_{s\varsigma}$ is the nonlinear partial contribution of the $r$th feature in the $\varsigma$th subspace. Similarly, extending Eq. 3, the generalized, nonlinear Incremental explanation is

$$\tilde{y}_i = \sum_r \left( f^{(r)}_{i0} + \sum_{\varsigma>0} [x \in s_{\varsigma}] \Delta f^{(r)}_{\Delta i \varsigma} \right) \qquad (7)$$

where $f_{i0}$ is the base contributions of the typical explanation model, and $\Delta f_{\Delta i \varsigma}$ is the incremental contributions of the $\varsigma$th special subspace explanation model. Here we consider that the incremental difference is additive, i.e., linear. Nonlinear effects could be investigated with multiplicative interactions ($x^2 \rightarrow x^3$) or a kernel transformation. To keep Incremental explanations simple, we can constrain the incremental contributions $\Delta f_{\Delta i \varsigma}$ with a sparsity regularization to reduce the number of terms, and with a smoothness regularization [2] to penalize overly curvy lines.

Next, we discuss generalizing Incremental explanations to models with interaction effects, i.e., multivariate functions that involve multiple features, e.g., $f(x^{(1)}, x^{(2)})$. Common models are rules and decision trees. While we have discussed several works to model subspaces with rules and trees, they do not support an incremental approach [30, 46]. To do so, future work could first convert any rule representation into a decision tree, compute the similarity between trees in each subspace (e.g., by calculating a graph edit distance [19]), and minimizing the difference. However, note that rules may overlap and lead to overlapping subspaces [29].

## 5.3 Scope of incremental explanations

We evaluated Incremental explanations for the understanding and memorability of explanatory factors in AI, specifically, for users to estimate the predictions that an AI would make (forward simulatability). This is meant to help human *cognition* toward decision making. Further work is needed to investigate whether these lead to improvements in downstream decision making, e.g., to decide whether to accept or reject a case based on quality estimations [64]; such a study would require careful framing and incentivization to ensure that users are correctly aligned and properly motivated to the task, and avoiding the confounder of prior knowledge which can diminish the benefits of XAI [38]. We do not propose it for *perception* tasks (e.g., vision and audio) or language reasoning (NLP), since they involve innate mental processes due to stimuli or low-level skills rather than deliberate reasoning.

Our paradigm of explanation incrementation assumes that users are novices who start with limited knowledge of the domain or

AI application, thus they need to be taught gently. We do not expect Incremental explanations to be strongly beneficial for domain experts who can handle complex data and have established conventions [34].

Similar to Poursabzi-Sangdeh et al. [48], we had evaluated Incremental explanations only for one application task of predicting housing prices. Perhaps, for applications that are less common (e.g., health diagnosis), or with critical but complicated numbers (e.g., decimals or fractions), Incremental and Subglobal explanations may still be overwhelming. Future work should validate our results across other application tasks.

## 5.4 Implications of incremental explanations

Our approach for Incremental explanations enables better learning of sparse linear factor models. This adds to the body of work of subspace-based explanations that take a divide-and-conquer approach to partitioning instances, and explaining each subspace as similarly as possible. COGAM moderated the number of visual chunks in line graphs [2], and it can be made to incrementally allow more curviness to allow users to learn more details. GlocalX [56] provides rule explanations in detail and in aggregate by merging them. Future work could investigate which explanation format (factors, line segments, or rules) are easier and more beneficial to learn incrementally.

Although our participants could well learn and recall Incremental explanations in our study, we acknowledge that the learning time was brief. Most learning that people do occurs over longer time periods with more repetitions. Hence, future work could deploy Incremental explanations to investigate its longitudinal benefits. We note that under longer durations, the learning of Subglobal explanations may also be improved, but perhaps less so than Incremental explanations, due to the slightly higher cognitive load.

Our approach of Incremental explanations limited the explanations to the same type (sparse linear models). However, users have diverse preferences for and usage strategies of explanations [36, 37], so incremented explanations should also be diverse. For example, first provide factors, then rules. This provides users with diverse retrieval cues, which can reinforce their memory of the explanations. Future work can explore how to increment across explanation structures.

## 6 CONCLUSION

We have introduced Incremental XAI to help users better recall and apply explanations of AI. This provides a set of base general factors for typical instances and sparse incremented factors for special cases. In modeling and user studies, we found that Incremental explanations help facilitate fast understanding like Global explanations that explain generally, and are easy to recall and understand like Subglobal explanations that explain subspaces more faithfully than Global explanations. Incremental explanations are also more memorable than Local explanations, facilitating better recall and understanding performance. This work demonstrates the importance of supporting more memorable explanations to deepen user understanding of AI for more productive interactions.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y Lim, and Mohan Kankanhalli. 2018. Trends and trajectories for explainable, accountable and intelligible systems: An hci research agenda. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–18.

[2] Ashraf Abdul, Christian von der Weth, Mohan Kankanhalli, and Brian Y Lim. 2020. COGAM: Measuring and Moderating Cognitive Load in Machine Learning Model Explanations. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–14.

[3] Amina Adadi and Mohammed Berrada. 2018. Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE access* 6 (2018), 52138–52160.

[4] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information fusion* 58 (2020), 82–115.

[5] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* 10, 7 (2015), e0130140.

[6] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S Lasecki, Daniel S Weld, and Eric Horvitz. 2019. Beyond accuracy: The role of mental models in human-AI team performance. In *Proceedings of the AAAI conference on human computation and crowdsourcing*, Vol. 7. 2–11.

[7] Anton Björklund, Andreas Henelius, Emilia Oikarinen, Kimmo Kallonen, and Kai Puolamäki. 2019. Sparse robust regression for explaining classifiers. In *Discovery Science: 22nd International Conference, DS 2019, Split, Croatia, October 28–30, 2019, Proceedings 22*. Springer, 351–366.

[8] Anton Björklund, Jarmo Mäkelä, and Kai Puolamäki. 2023. SLISEMAP: Supervised dimensionality reduction through local explanations. *Machine Learning* 112, 1 (2023), 1–43.

[9] Leo Breiman. 2001. Random forests. *Machine learning* 45 (2001), 5–32.

[10] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–21.

[11] Ruth MJ Byrne. 2019. Counterfactuals in Explainable Artificial Intelligence (XAI): Evidence from Human Reasoning.. In *IJCAI*. 6276–6282.

[12] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. 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*. 1721–1730.

[13] Marco Cerliani. 2022. linear-trees. https://github.com/cerlymarco/linear-tree.

[14] Tanya Chowdhury, Razieh Rahimi, and James Allan. 2022. Equi-explanation maps: concise and informative global summary explanations. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 464–472.

[15] Anupam Datta, Shayak Sen, and Yair Zick. 2016. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In *2016 IEEE symposium on security and privacy (SP)*. IEEE, 598–617.

[16] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).

[17] Pierre Dragicevic. 2018. Can we call mean differences "effect sizes". https://transparentstatistics.org/2018/07/05/meanings-effect-size. Accessed: 2023-12-10.

[18] Jeff Druce, James Niehaus, Vanessa Moody, David Jensen, and Michael L Littman. 2021. Brittle AI, causal confusion, and bad mental models: challenges and successes in the XAI program. *arXiv preprint arXiv:2106.05506* (2021).

[19] Xinbo Gao, Bing Xiao, Dacheng Tao, and Xuelong Li. 2010. A survey of graph edit distance. *Pattern Analysis and applications* 13 (2010), 113–129.

[20] Shirley Gregor and Izak Benbasat. 1999. Explanations from intelligent systems: Theoretical foundations and implications for practice. *MIS quarterly* (1999), 497–530.

[21] Harlfoxem. 2016. House Sales in King County, USA. https://www.kaggle.com/harlfoxem/housesalesprediction.

[22] Fritz Heider. 2013. *The psychology of interpersonal relations*. Psychology Press.

[23] Andras Janosi, William Steinbrunn, Matthias Pfisterer, and Robert Detrano. 1988. Heart Disease. UCI Machine Learning Repository. DOI: https://doi.org/10.24432/C52P4X.

[24] Harmanpreet Kaur, Eytan Adar, Eric Gilbert, and Cliff Lampe. 2022. Sensible AI: Re-imagining Interpretability and Explainability using Sensemaking Theory. In *2022 ACM Conference on Fairness, Accountability, and Transparency*.

[25] Sunnie SY Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, and Andrés Monroy-Hernández. 2023. " Help Me Help the AI": Understanding How Explainability Can Support Human-AI Interaction. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.

[26] Rafal Kocielnik, Saleema Amershi, and Paul N Bennett. 2019. Will you accept an imperfect ai? exploring designs for adjusting end-user expectations of ai systems. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–14.

[27] Josua Krause, Adam Perer, and Kenney Ng. 2016. Interacting with predictions: Visual inspection of black-box machine learning models. In *Proceedings of the 2016 CHI conference on human factors in computing systems*. 5686–5697.

[28] Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Samuel J Gershman, and Finale Doshi-Velez. 2019. Human evaluation of models built for interpretability. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 1. 59–67.

[29] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. 2016. Interpretable decision sets: A joint framework for description and prediction. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1675–1684.

[30] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. 2019. Faithful and customizable explanations of black box models. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 131–138.

[31] Markus Langer, Daniel Oster, Timo Speith, Holger Hermanns, Lena Kästner, Eva Schmidt, Andreas Sesing, and Kevin Baum. 2021. What do we want from Explainable Artificial Intelligence (XAI)?–A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research. *Artificial Intelligence* 296 (2021), 103473.

[32] Benjamin Letham, Cynthia Rudin, Tyler H McCormick, and David Madigan. 2015. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. (2015).

[33] Q Vera Liao and Kush R Varshney. 2021. Human-centered explainable ai (xai): From algorithms to user experiences. *arXiv preprint arXiv:2110.10790* (2021).

[34] Brian Y Lim, Joseph P Cahaly, Chester YF Sng, and Adam Chew. 2023. Diagrammatization: Rationalizing with diagrammatic AI explanations for abductive reasoning on hypotheses. *arXiv preprint arXiv:2302.01241* (2023).

[35] Brian Y Lim and Anind K Dey. 2009. Assessing demand for intelligibility in context-aware applications. In *Proceedings of the 11th international conference on Ubiquitous computing*. 195–204.

[36] Brian Y Lim and Anind K Dey. 2011. Design of an intelligible mobile context-aware application. In *Proceedings of the 13th international conference on human computer interaction with mobile devices and services*. 157–166.

[37] Brian Y Lim and Anind K Dey. 2013. Evaluating intelligibility usage and usefulness in a context-aware application. In *International Conference on Human-Computer Interaction*. Springer, 92–101.

[38] Brian Y Lim, Anind K Dey, and Daniel Avrahami. 2009. Why and why not explanations improve the intelligibility of context-aware intelligent systems. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 2119–2128.

[39] Tania Lombrozo. 2006. The structure and function of explanations. *Trends in cognitive sciences* 10, 10 (2006), 464–470.

[40] Tania Lombrozo. 2007. Simplicity and probability in causal explanation. *Cognitive psychology* 55, 3 (2007), 232–257.

[41] Ana Lucic, Hinda Haned, and Maarten de Rijke. 2020. Why does my model fail? contrastive local explanations for retail forecasting. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 90–98.

[42] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems* 30 (2017).

[43] Yan Lyu, Fan Gao, I-Shuen Wu, and Brian Y Lim. 2020. Imma Sort by Two or More Attributes With Interpretable Monotonic Multi-Attribute Sorting. *IEEE Transactions on Visualization and Computer Graphics* 27, 4 (2020), 2369–2384.

[44] Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence* 267 (2019), 1–38.

[45] Menaka Narayanan, Emily Chen, Jeffrey He, Been Kim, Sam Gershman, and Finale Doshi-Velez. 2018. How do humans understand explanations from machine learning systems? an evaluation of the human-interpretability of explanation. *arXiv preprint arXiv:1802.00682* (2018).

[46] Karthikeyan Natesan Ramamurthy, Bhanukiran Vinzamuri, Yunfeng Zhang, and Amit Dhurandhar. 2020. Model agnostic multilevel explanations. *Advances in neural information processing systems* 33 (2020), 5968–5979.

[47] Mahsan Nourani, Chiradeep Roy, Jeremy E Block, Donald R Honeycutt, Tahrima Rahman, Eric D Ragan, and Vibhav Gogate. 2022. On the Importance of User Backgrounds and Impressions: Lessons Learned from Interactive AI Applications. *ACM Transactions on Interactive Intelligent Systems* 12, 4 (2022), 1–29.

[48] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Wortman Vaughan, and Hanna Wallach. 2021. Manipulating and measuring model interpretability. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–52.

[49] J. Ross Quinlan. 1986. Induction of decision trees. *Machine learning* 1 (1986), 81–106.

[50] R. Quinlan. 1993. Auto MPG. UCI Machine Learning Repository. DOI: https://doi.org/10.24432/C5859H.

[51] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.

[52] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.

[53] Nathalie Henry Riche, Christophe Hurter, Nicholas Diakopoulos, and Sheelagh Carpendale. 2018. *Data-driven storytelling*. CRC Press.

[54] Jason D Rights and Sonya K Sterba. 2019. Quantifying explained variance in multilevel models: An integrative framework for defining R-squared measures. *Psychological methods* 24, 3 (2019), 309.

[55] Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. 2022. Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys* 16 (2022), 1–85.

[56] Mattia Setzu, Riccardo Guidotti, Anna Monreale, Franco Turini, Dino Pedreschi, and Fosca Giannotti. 2021. Glocalx-from local to global explanations of black box ai models. *Artificial Intelligence* 294 (2021), 103457.

[57] Thomas J Shuell. 1986. Cognitive conceptions of learning. *Review of educational research* 56, 4 (1986), 411–436.

[58] Aaron Springer and Steve Whittaker. 2020. Progressive disclosure: When, why, and how do users want algorithmic transparency information? *ACM Transactions on Interactive Intelligent Systems (TiiS)* 10, 4 (2020), 1–32.

[59] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International conference on machine learning*. PMLR, 3319–3328.

[60] Chun-Hua Tsai, Yue You, Xinning Gui, Yubo Kou, and John M Carroll. 2021. Exploring and promoting diagnostic transparency and explainability in online symptom checkers. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–17.

[61] Berk Ustun and Cynthia Rudin. 2016. Supersparse linear integer models for optimized medical scoring systems. *Machine Learning* 102 (2016), 349–391.

[62] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* 31 (2017), 841.

[63] Danding Wang, Qian Yang, Ashraf Abdul, and Brian Y Lim. 2019. Designing theory-driven user-centric explainable AI. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–15.

[64] Danding Wang, Wencan Zhang, and Brian Y Lim. 2021. Show or suppress? Managing input uncertainty in machine learning model explanations. *Artificial Intelligence* 294 (2021), 103456.

[65] Xinru Wang and Ming Yin. 2023. Watch Out for Updates: Understanding the Effects of Model Explanation Updates in AI-Assisted Decision Making. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–19.

[66] Yunlong Wang, Priyadarshini Venkatesh, and Brian Y Lim. 2022. Interpretable Directed Diversity: Leveraging Model Explanations for Iterative Crowd Ideation. In *CHI Conference on Human Factors in Computing Systems*. 1–28.

[67] Ming Yin, Jennifer Wortman Vaughan, and Hanna Wallach. 2019. Understanding the effect of accuracy on trust in machine learning models. In *Proceedings of the 2019 chi conference on human factors in computing systems*. 1–12.

[68] Wencan Zhang and Brian Y Lim. 2022. Towards Relatable Explainable AI with the Perceptual Process. In *CHI Conference on Human Factors in Computing Systems*. 1–24.

# A   APPENDIX

## A.1   Decision surfaces of each explanation

In A.1, we examine the decision surfaces of the predictor and explainer models to study how the different XAI types support linear or piecewise-linear relationships between multiple attributes and the prediction value. Based on the House Sales dataset, this provides a conceptual interpretation to the reader of i) how nonlinear the decision surface of the AI Predictor model is (a), and how the Local explanations cumulatively capture the nonlinearity (e); ii) how Subspace explanations (c) better fit the nonlinear decision surface of the AI predictor model (e) compared to the Global explanation (b); and iii) how Incremental explanations simplify memorability by first conveying a linear relationship (d, left), then showing a partial linear segment (d, right)
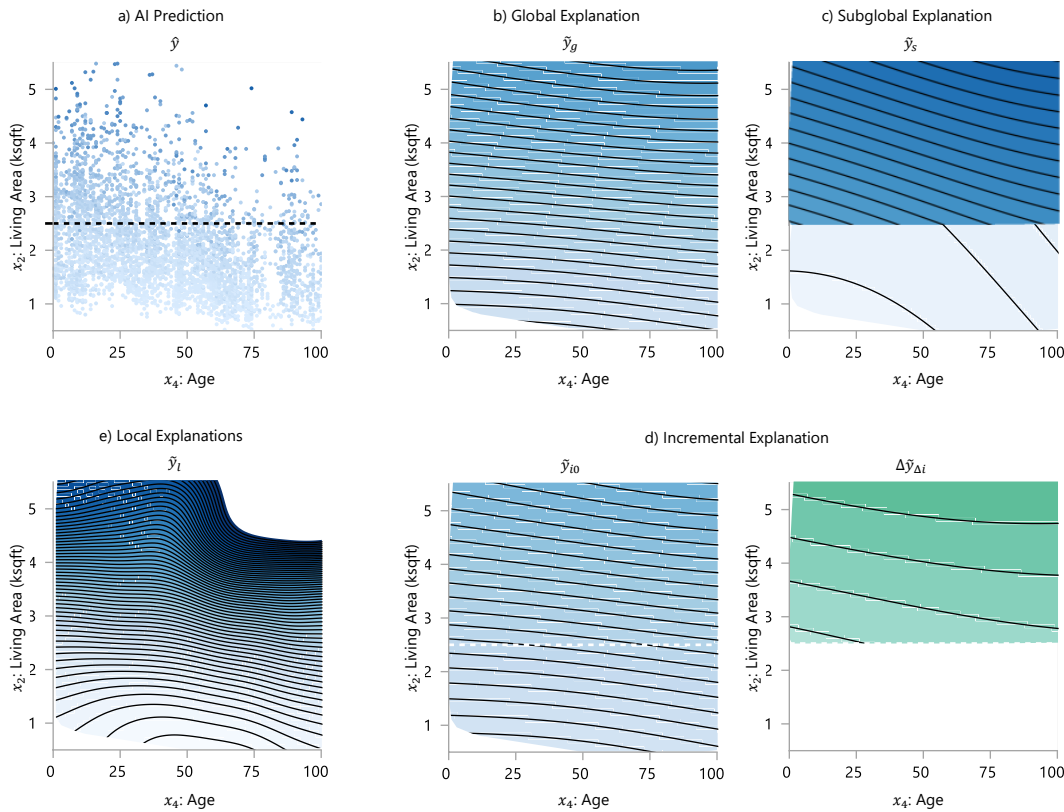


**Figure A.1: Decision surface of AI prediction model and various explanation models with two attributes (x-y axes), showing prediction output (color: darker is higher value) for the House Sales in King County, USA [21] dataset. Models are based on the dataset used in the modeling and user studies; showing only values for attributes 2 and 4. a) Scatter plot of AI System predictions for instances with various $(x_2, x_4)$ values. Dashed line indicates threshold to split for Subglobal and Incremental explanation models. b) Contour plot of Global explanation showing linear slope mostly increasing in the direction of $x_2$ i.e., $w_2 > w_4$ (see factors in Fig. 2). c) Contour plot of Subglobal explanation model showing two linear models – gentler-sloped for typical instances ($x_2 < 2.5$) and steeper-sloped for outlier instances ($x_2 \geq 2.5$). d) Contour plot of Incremental explanation model showing general model for all instances (Left), and incrementally-sloped model for outlier instances ($x_2 \geq 2.5$, Right). e) Contour plot of accumulation of multiple Local explanation models showing non-linear surface that manifests when learning from heterogeneous local explanations.**

## A.2 Subglobal and Incremental explanations performance across varying subspace thresholds

Subglobal and Incremental explanations model data with two aspects: first partitioning the feature space into subspaces, then fitting linear factors for each subspace. Although the partition thresholds are learned automatically with the linear model tree, they may be adjusted to simplify presentation to users. We examined how selecting different partition thresholds affect the faithfulness of each subspace explainer model, how the factors change, and whether incremental factors are kept small.



Figure A.2: Results of XAI faithfulness and factors for Subglobal and Incremental explanations when varying subspace partitioning thresholds, for different datasets and subspaces. In the top row graphs, lower MAE is better, and normalized MAE is shown to clarify the curve minima. In the bottom graphs, incremental factors are regularized to 0 if the weight values are the same for subspace $s_<$ as for $s_\geq$. Vertical black lines indicate the optimal partition threshold chosen in the modeling and user studies. Gray horizontal lines indicate the factors for the Global explanation.

For this experiment, we manually set the partition threshold of the feature selected by the linear model tree (Living Area $x^{(2)}$ for House Sales, Age $x^{(1)}$ for Heart Disease, Horsepower $x^{(3)}$ for Auto MPG), then trained Subglobal and Incremental models for each threshold level. We varied thresholds from the 10th to 90th percentile of feature values, see Fig. 6 for results. We denote the subspace below the threshold as $s_<$ and the other as $s_\geq$. Note that as the threshold changes, the ratio of subspace sizes changes, so we do not label them as typical or outlier. We calculated the Combined MAE as the sample-weighted sum of subspace MAE based on the number of instances in each subspace, then normalize each MAE metric for clarity (see Fig. 6, top). The lowest Combined MAE coincides with the thresholds selected by the linear model tree. Fig. 6 (bottom) shows how each linear factor $w^{(r)}$ changes with partition threshold.

As expected, the weights across XAI type for each subspace is somewhat similar (i.e., $w_{s\varsigma}^{(r)} \approx w_{i\varsigma}^{(r)}$), with differences due to different loss functions. However, the weights of Subglobal explanations across subspaces are generally different (i.e., $w_{s<}^{(r)} \neq w_{s\geq}^{(r)}$), since each subspace is best fit with different linear models. Conversely, the weights of Incremental explanations are sometimes similar across subspaces (i.e., $\Delta w_i^{(r)} = w_{i<}^{(r)} - w_{i\geq}^{(r)} = 0$), indicating that some incremental factors have been regularized to 0. This is seen when the dark blue and red lines overlap in Fig. 6 (bottom). Note how the MAE may not be the lowest for such cases, so there is a trade-off to prioritize performance over simplicity (fewer Incremental factors), which is typical in XAI methods [4] For House Sales, the incremental factors $\Delta w_i^{(1)}$ and $\Delta w_i^{(3)}$ are always almost 0, and $\Delta w_i^{(2)} = 0$ when Living Area < 1.8 ksqft. In general, Base+Incremental factors are not identical to Global or Subglobal factors due to different training objectives.

## A.3   Qualtrics survey for the summative user study

Figures A.3-A.21 depict the Qualtrics survey workflow with each of the four XAI conditions for the summative user study.



**Figure A.3: Introduction of housing price estimation.**



**Figure A.4: Introduction of attributes in housing price estimation.**

## Prior Experience

Please answer the following questions on your prior experience with housing price estimation and AI technologies. There are no right or wrong answers, so be as truthful as possible.

I would rate my experience in <u>housing price estimation</u> as:
◯ I have never attempted the task before and have little or no knowledge of it.
◯ I am familiar with the task and have some understanding on how to perform it.
◯ I am proficient in the task and perform it somewhat often.
◯ I have professional working knowledge in the task.

I would rate my experience with <u>AI technologies</u> as:
◯ I have no formal education/training in AI or related topics, and I do not have much knowledge of recent trends in AI.
◯ I have no formal education/training in AI or related topics, but I do have a knowledge of recent trends in AI.
◯ I have no formal education/training in AI, but I have training in related technical topics (e.g., statistics, math, software engineering, technology).
◯ I have some education/training specific to AI.

I would rate my interest in <u>AI technologies</u> as:
◯ I have no interest in understanding or using AI.
◯ I have interest in understanding AI but not using it.
◯ I have interest in using AI but not understanding it.
◯ I have interest in both understanding and using AI.

**Figure A.5: Questions on users' prior experience with housing price estimation and AI background.**

## Tutorial: AI Explainer

The AI Explainer shows you how each attribute contributes to the AI System's predicted price. It estimates a <u>factor</u> for each attribute and a price <u>adjustment</u>.

**Important Note**: For simplicity, we <u>round</u> the numbers in the AI Explainer interface, but use precise numbers in calculations. You can see the exact numbers by hovering your mouse over each number.

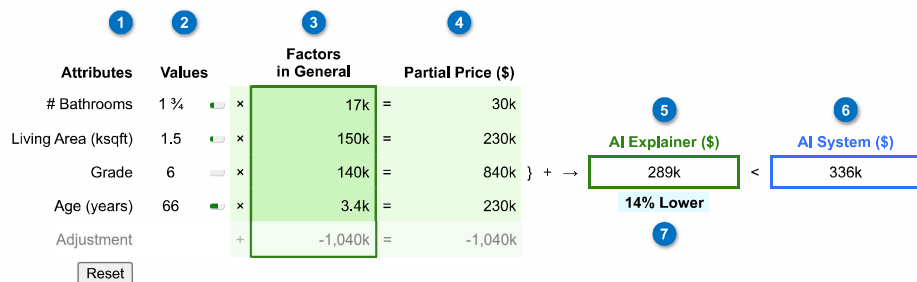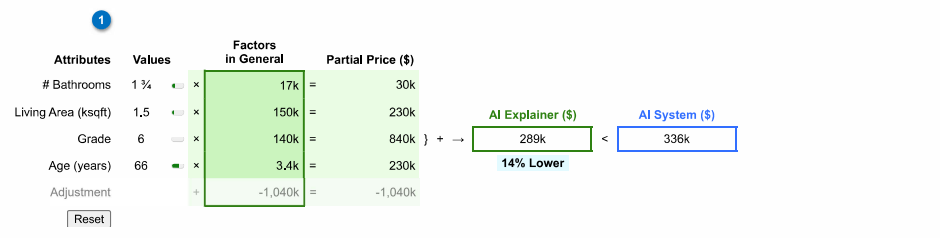CLICK NEXT to see components of the basic AI Explainer in detail.



**Figure A.6: Tutorial on basic (Global) AI explanation.**

## Tutorial: AI Explainer

### Basic AI Explainer: Attributes and Values
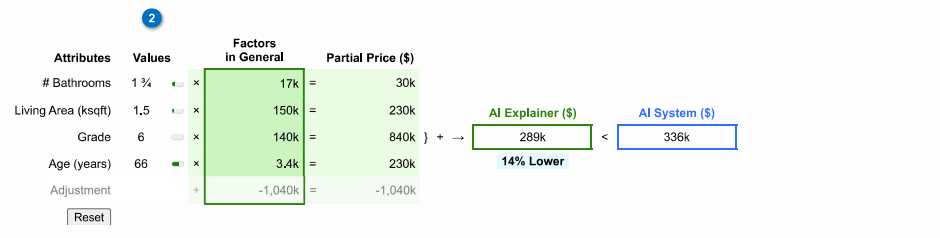
**1.** The list of the attributes that you need to know. The Adjustment is an additional value calculated based on other unlisted attributes.



**2.** The value of each attribute. The green bar beside each attribute value shows how *big or small* it is relative other houses.

For example, this house:

- has 1 ¾ bathrooms, which is lower than average (short green bar)
- is 66 years old, which is older than average (long green bar)



## Tutorial: AI Explainer

### Basic AI Explainer: Factors and Partial Prices

**3.** The factors to multiply with each attribute.

- 1 bathroom is worth $17k, each thousand sqft of area is worth $150k, etc.



**4.** The partial price contributed by each attribute, calculated as *Value × Factor*.

- 1 ¾ bathrooms contribute $30k, the age of 66 years contributes 230k, etc.
- The adjustment factor is −$1,040k. which is negative so this is subtracted from the price.



**Figure A.7: Tutorial on the attributes, values, factors, and partial prices of the explanation interface. All participants are first trained on the Global explanation since it has the simplest format.**

# Tutorial: AI Explainer

## Basic AI Explainer: Estimated and True Prices

**5.** The AI Explainer estimated price, calculated by summing all the partial prices and adjustment.

**Reminder**: For simplicity, we round the numbers in the AI Explainer interface, but use precise numbers in calculations. You may see the exact numbers by hovering your mouse over each number.

$30k + $230k + $840k + $230k − $1,040k = $289k



**6.** The correct price predicted by the AI System.



**7.** The difference between the AI Explainer and the AI System. The AI Explainer approximates a simplified version of AI System and thus has an error.

Here, it estimates 14% lower than the correct value.



**Figure A.8: Tutorial on estimated explanation price, correct AI System price, and the percent difference.**

## Screening Test

To continue to the rest of the survey, we need to make sure you understand how to use the AI Explainer. Please correctly answer these screening questions.

| Attributes | Values | | Factors in General | | Partial Price ($) |
|---|---|---|---|---|---|
| # Bathrooms | 1 ¾ | × | 17k | = | 30k |
| Living Area (ksqft) | 1.5 | × | 150k | = | 230k |
| Grade | 6 | × | 140k | = | 840k |
| Age (years) | 66 | × | 3.4k | = | 230k |
| Adjustment | | + | -1,040k | = | -1,040k |

AI Explainer ($): 289k  <  AI System ($): 336k

**14% Lower**

} + →

Reset

Please answer the following questions correctly to continue.

What is the <u>Value</u> of Living Area?

◯ 1   ◯ 150   ◯ 1.5   ◯ 230   ◯ 6

What is the <u>Factor</u> for # Bathrooms?

◯ 0k   ◯ 150k   ◯ 17k   ◯ 840k   ◯ 140k

What is the <u>Partial Price</u> for # Grade?

◯ 0k   ◯ 150k   ◯ 17k   ◯ 840k   ◯ 140k

What is the output of the AI Explainer here. If your answer is within **10%**, it will be accepted.

| 0k | 200k | 400k | 600k | 800k | 1,000k | 1,200k | 1,400k | 1,600k | 1,800k | 2,000k |

What is the output of the AI System here. If your answer is within **10%**, it will be accepted.

| 0k | 200k | 400k | 600k | 800k | 1,000k | 1,200k | 1,400k | 1,600k | 1,800k | 2,000k |

**Figure A.9: Screening questions for the Global explanation to check users' comprehension.**

# Tutorial: AI Explainer

## Full AI Explainer

The AI Explainer can be made <u>more accurate</u> by showing multiple sets of factors instead of one.

Here, the basic AI Explainer adds the rules "**A**" and "**B**". The factors are different depending on rule condition. Each house belongs to exactly one rule condition.

- If **Living Area < 2.5 ksqft** (*top example where Living Area = 1.5*), use the factors in Explanation **A**.
- Otherwise, if **Living Area ≥ 2.5 ksqft** (*bottom example where Living Area = 2.8*), use the factors in Explanation **B**.

Notice how this is closer to the AI System's prediction than the basic AI Explainer.

Rule A:

| Attributes | Values | | | Factors when Living Area < 2.5 | | Partial Price ($) | | | | AI Explainer ($) | | AI System ($) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 1 ¾ | | × | 16k | = | 28k | | | | | | |
| Living Area (ksqft) | 1.5 | | × | 63k | = | 96k | | | | | | |
| Grade | 6 | | × | 120k | = | 720k | } | + | → | 314k | < | 336k |
| Age (years) | 66 | | × | 2.5k | = | 160k | | | | **6.6% Lower** | | |
| Adjustment | | | + | -697k | = | -697k | | | | | | |

Reset

Rule B:

| Attributes | Values | | | Factors when Living Area ≥ 2.5 | | Partial Price ($) | | | | AI Explainer ($) | | AI System ($) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 2 ½ | | × | 57k | = | 140k | | | | | | |
| Living Area (ksqft) | 2.8 | | × | 240k | = | 670k | | | | AI Explainer ($) | | AI System ($) |
| Grade | 9 | | × | 160k | = | 1,400k | } | + | → | 661k | > | 645k |
| Age (years) | 16 | | × | 4.9k | = | 78k | | | | **2.4% Higher** | | |
| Adjustment | | | + | -1,660k | = | -1,660k | | | | | | |

Reset

**Figure A.10: Tutorial on the Subglobal explanation.**

## Screening Test

To continue to the rest of the survey, we need to make sure you understand how to use the AI Explainer. Please correctly answer these screening questions.

This is an example following Rule A:

| Attributes | Values | | Factors when Living Area < 2.5 | | Partial Price ($) |
|---|---|---|---|---|---|
| # Bathrooms | 1 ¾ | × | 16k | = | 28k |
| Living Area (ksqft) | 1.5 | × | 63k | = | 96k |
| Grade | 6 | × | 120k | = | 720k |
| Age (years) | 66 | × | 2.5k | = | 160k |
| Adjustment | | + | -697k | = | -697k |

Reset

**AI Explainer ($)**  314k  <  **AI System ($)**  336k
720k } + →
6.6% Lower

What is the Value of Living Area?

○ 66  ○ 1.8  ○ 7  ○ 6  ○ 1.5

What is the Factor when Living Area < 2.5 for # Bathrooms?

○ 720k  ○ 2.5k  ○ 16k  ○ 36k  ○ 120k

What is the Partial Price for Grade?

○ 720k  ○ 2.5k  ○ 16k  ○ 36k  ○ 120k

Based on the attribute of the following four houses, which rule should the AI Explainer use for each one?

| | A | B | Not enough information |
|---|---|---|---|
| 1) Living Area = 1.8 ksqft | ○ | ○ | ○ |
| 2) Living Area = 2.5 ksqft | ○ | ○ | ○ |
| 3) Living Area = 3.1 ksqft | ○ | ○ | ○ |
| 4) Grade = 11 | ○ | ○ | ○ |

What is the output of the AI Explainer. If your answer is within **10%**, it will be accepted.

Ok  200k  400k  600k  800k  1,000k  1,200k  1,400k  1,600k  1,800k  2,000k

What is the output of the AI System. If your answer is within **10%**, it will be accepted.

Ok  200k  400k  600k  800k  1,000k  1,200k  1,400k  1,600k  1,800k  2,000k

**Figure A.11: Screening questions for the Subglobal explanation to check users' comprehension.**
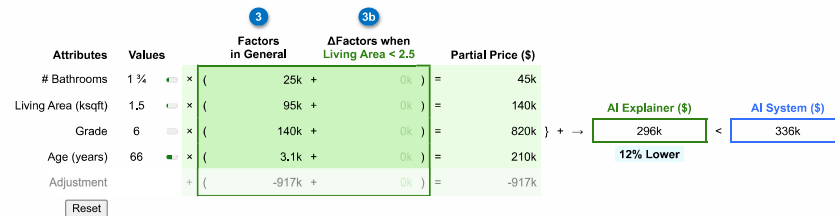
## Tutorial: AI Explainer

### Full AI Explainer: General Houses

The AI Explainer can be made more accurate by showing modified factors for special types of houses.

(3) The first column has the **General** factors, which always remains the same for all houses.
(3b) The second column has the factor adjustments, which are added to the **General** factors if and only if the **Special Case** rule applies.

Here, AI Explainer modifies factors for houses defined by the **Special Case** rule of **Living Area ≥ 2.5 ksqft**. In this example, since Living Area = 1.5 ksqft, the **Special Case** does not apply and the adjustments are 0k. This is essentially the same as the basic AI Explainer shown before.

The next page shows a house for which the **Special Case** rule applies to.

| Attributes | Values | | | Factors in General (3) | | ΔFactors when Living Area < 2.5 (3b) | | Partial Price ($) |
|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 1 ¾ | | × ( | 25k | + | 0k | ) = | 45k |
| Living Area (ksqft) | 1.5 | | × ( | 95k | + | 0k | ) = | 140k |
| Grade | 6 | | × ( | 140k | + | 0k | ) = | 820k |
| Age (years) | 66 | | × ( | 3.1k | + | 0k | ) = | 210k |
| Adjustment | | | + ( | -917k | + | 0k | ) = | -917k |

**AI Explainer ($)**: 296k  <  **AI System ($)**: 336k
**12% Lower**

[Reset]

## Tutorial: AI Explainer

### Full AI Explainer: Special Case Houses

(3) Notice that the **General** factors remain the same as the previous example.
(3b) However, since Living Area = 5.9 ksqft, the **Special Case** factor adjustments are added to the **General** factors. Not all attributes require adjustments.

*For example, the final factor for Living Area is $95k + $120 = $215k, which means that each ksqft of area is worth more for a large house than a small house.*

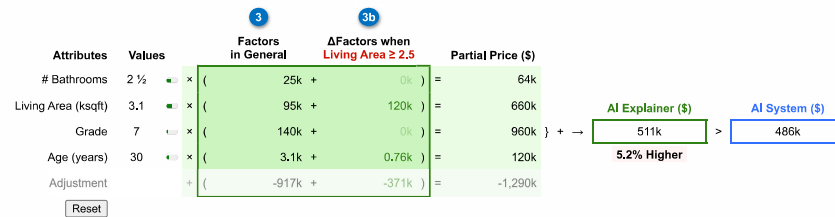Notice how this is closer to the AI System's prediction than the basic AI Explainer.

| Attributes | Values | | | Factors in General (3) | | ΔFactors when Living Area ≥ 2.5 (3b) | | Partial Price ($) |
|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 2 ½ | | × ( | 25k | + | 0k | ) = | 64k |
| Living Area (ksqft) | 3.1 | | × ( | 95k | + | 120k | ) = | 660k |
| Grade | 7 | | × ( | 140k | + | 0k | ) = | 960k |
| Age (years) | 30 | | × ( | 3.1k | + | 0.76k | ) = | 120k |
| Adjustment | | | + ( | -917k | + | -371k | ) = | -1,290k |

**AI Explainer ($)**: 511k  >  **AI System ($)**: 486k
**5.2% Higher**

[Reset]

**Figure A.12: Tutorial on the Incremental explanation.**

## Screening Test

To continue to the rest of the survey, we need to make sure you understand how to use the AI Explainer. Please correctly answer these screening questions.



What is the Value of Living Area?

○ 0   ○ 95   ○ 2.8   ○ 120   ○ 2.5

What is the Factor in General for # Bathrooms?

○ 0k   ○ 140k   ○ 25k   ○ 296k   ○ 95k

What is the ΔFactor Adjustment when Living Area ≥ 2.5 for # Bathrooms?

○ 0k   ○ 180k   ○ 21k   ○ 3k   ○ 40k

What is the Partial Price for Grade?

○ 16k   ○ 1,200k   ○ 9k   ○ 64k   ○ 40k

Based on the attribute of the following four houses, which AI Explainer should be used for each one?

| | General | Special Case | Not enough information |
|---|---|---|---|
| 1) Living Area = 1.8 ksqft | ○ | ○ | ○ |
| 2) Living Area = 2.5 ksqft | ○ | ○ | ○ |
| 3) Living Area = 3.1 ksqft | ○ | ○ | ○ |
| 4) Grade = 11 | ○ | ○ | ○ |

What is the output of the AI Explainer If your answer is within **10%**, it will be accepted.

0k  200k  400k  600k  800k  1,000k  1,200k  1,400k  1,600k  1,800k  2,000k

**Figure A.13: Screening questions for the Incremental explanation to check users' comprehension.**

## Tutorial: AI Explainer

### Full AI Explainer

The AI Explainer becomes more accurate when its factors change based on the attribute values. Here are some example houses with different attributes and factors.

**Helpful hints**:

- Houses with similar attribute values should have similar factors.
- The different factor levels affects the partial price contributions of each attribute (e.g., Grade is worth much less in Examples 1 and 2 than in Example 3).
- Negative factors mean that the attribute lowers the price of that particular house.

Example 1: Small, old house with a poor grade.

| Attributes | Values | | | Factors when Similar Case | | Partial Price ($) | | AI Explainer ($) | | AI System ($) |
|---|---|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 1 ¾ | | × | -2.2k | = | -3.8k | | | | |
| Living Area (ksqft) | 1.5 | | × | -18k | = | -27k | | 336k | < | 336k |
| Grade | 6 | | × | -7.1k | = | -43k | } + → | | | |
| Age (years) | 66 | | × | 0.19k | = | 12k | | 0.077% Lower | | |
| Adjustment | | | + | 396k | = | 396k | | | | |

Example 2: Small, modern house with an average grade.

| Attributes | Values | | | Factors when Similar Case | | Partial Price ($) | | AI Explainer ($) | | AI System ($) |
|---|---|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 2 | | × | -1.2k | = | -2.5k | | | | |
| Living Area (ksqft) | 1.8 | | × | -8.5k | = | -16k | | 454k | > | 436k |
| Grade | 8 | | × | 2.3k | = | 19k | } + → | | | |
| Age (years) | 21 | | × | -0.29k | = | -6.2k | | 4% Higher | | |
| Adjustment | | | + | 459k | = | 459k | | | | |

Example 3: Large, new house with a high grade.

| Attributes | Values | | | Factors when Similar Case | | Partial Price ($) | | AI Explainer ($) | | AI System ($) |
|---|---|---|---|---|---|---|---|---|---|---|
| # Bathrooms | 3 | | × | -6.3k | = | -19k | | | | |
| Living Area (ksqft) | 5.9 | | × | 77k | = | 450k | | 1,870k | > | 1,700k |
| Grade | 11 | | × | 94k | = | 1,000k | } + → | | | |
| Age (years) | 11 | | × | 6.4k | = | 71k | | 9.5% Higher | | |
| Adjustment | | | + | 328k | = | 328k | | | | |

**Figure A.14: Tutorial on the Local explanation.**

## Screening Test

To continue to the rest of the survey, we need to make sure you understand how to use the AI Explainer. Please correctly answer these screening questions.



What is the <u>Value</u> of Living Area?

○ 1.5    ○ 396    ○ 1.75    ○ 336    ○ 6

What is the <u>Factor when Similar Case</u> for # Bathrooms?

○ −43k    ○ −3.8k    ○ −18k    ○ −2.2k    ○ 1.2k

What is the <u>Partial Price</u> for Grade?

○ −43k    ○ −3.8k    ○ −18k    ○ −2.2k    ○ 1.2k

What is the output of the AI Explainer. If your answer is within **10%**, it will be accepted.



What is the output of the AI System here. If your answer is within **10%**, it will be accepted.



**Figure A.15: Screening questions for the Local explanation to check users' comprehension.**

**Test Case 1**

**Examine this case that the AI System would predict on.**

**Hints:**

1. Try to recall and estimate the **factors** you memorized and write them into the textboxes in the AI Explainer.
2. The partial prices and the AI Explainer output are **calculated automatically**.
3. Use the calculation results to **estimate** what the AI System would have predicted.



**What do you think are the most influential attributes to determine the price in this case?**

| | Select your best guess: | | | | The attribute **value** is...<br>**Hint:** look at the green bar beside the value. | | | Does this attribute increase or decrease the price?<br>**Hint:** look at the **factor** sign. | | | By an estimate of how much $k per unit? |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Bathrooms | Living Area | Grade | Age | Small | Medium | Large | Increase (positive) | No Change | Decrease (negative) | Enter the factor (Number only): |
| Most influential: | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ▢ |
| 2nd most influential: | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ▢ |

**What do you think might be the AI Explainer's behaviour?**

| | I think the AI Explainer might be [fill in the blank] compared to the AI System. | | | I think the difference % between the AI Explainer the AI System might be: | | |
|---|---|---|---|---|---|---|
| | Lower | Similar | Higher | Small (< 5%) | Moderate (5–10%) | Large (> 10%) |
| Select: | ○ | ○ | ○ | ○ | ○ | ○ |

**What do you think is the output of the AI Explainer? If your answer is within 10%, it will be accepted.**

| 0k | 200k | 400k | 600k | 800k | 1,000k | 1,200k | 1,400k | 1,600k | 1,800k | 2,000k |
|---|---|---|---|---|---|---|---|---|---|---|

**What do you think is the output of the AI System prediction? If your answer is within 10%, it will be accepted.**

| 0k | 200k | 400k | 600k | 800k | 1,000k | 1,200k | 1,400k | 1,600k | 1,800k | 2,000k |
|---|---|---|---|---|---|---|---|---|---|---|

**Figure A.16: Sample of the unassisted forward simulation trial, where participants are asked to estimate the explanation and AI system outputs. The displayed UI is the Incremental condition.**

**Test Case 1**

Now study the AI Explanation to help you understand its prediction from the AI System.

| Attributes | Values | | | Factors in General | ΔFactors when Living Area ≥ 2.5 | Partial Price ($) |
|---|---|---|---|---|---|---|
| # Bathrooms | 2 ½ | × | ( | 25k + | 0k ) = | 64k |
| Living Area (ksqft) | 2.6 | × | ( | 95k + | 120k ) = | 570k |
| Grade | 8 | × | ( | 140k + | 0k ) = | 1,100k } + → |
| Age (years) | 16 | × | ( | 3.1k + | 0.76k ) = | 62k |
| Adjustment | | + | ( | -917k + | -371k ) = | -1,290k |

AI Explainer ($): 499k  ≈  AI System ($): ?

Reset

What are the <u>most influential</u> attributes to determine the price <u>in this case</u>?
**Hint**: the most influential attribute is the one with the largest partial price, either positive or negative.

| | Select the appropriate attribute: | | | | The attribute **value** is... Hint: look at the green bar beside the value. | | | Does this attribute increase or decrease the price? Hint: look at the **factor** sign. | | | By how much $k per unit? |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Bathrooms | Living Area | Grade | Age | Small | Medium | Large | Increase (positive) | No Change | Decrease (negative) | Enter the factor (Number only): |
| <u>Most</u> influential: | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | [ ] |
| <u>2nd most</u> influential: | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | [ ] |

What do you think might be the **AI Explainer**'s behaviour?

| | I think the AI Explainer might be [fill in the blank] compared to the AI System. | | | I think the difference % between the AI Explainer the AI System might be: | | |
|---|---|---|---|---|---|---|
| | Lower | Similar | Higher | Small (< 5%) | Moderate (5–10%) | Large (> 10%) |
| Select: | ○ | ○ | ○ | ○ | ○ | ○ |

What do you think is the output of the **AI System** prediction? If your answer is within 10%, it will be accepted.

| 0k | 200k | 400k | 600k | 800k | 1,000k | 1,200k | 1,400k | 1,600k | 1,800k | 2,000k |
|---|---|---|---|---|---|---|---|---|---|---|

**Figure A.17: Sample of the assisted forward simulation trial, where participants are asked to estimate the AI system output based on the given explanation. The displayed UI is the Incremental condition.**

## Check your answer

**Here is what the AI System predicted. Please study and compare your answers with the correct ones.**



Notice that each **factor** is the **same** for each case in the **same group of cases** (typical **or** exceptional). Notice how **some factors** (which?) are the **same** for **exceptional situations**.

**Real AI Explainer estimation**: **$527k**.
Your AI Explainer prediction: **$828k**
The difference is: **+$301k**, 57.2% Higher.
Estimation is <u>wrongly outside 10%</u> of the error margin.

**Real AI System prediction**: **$511k**.
Your AI System prediction **without AI Explainer**: **$681k**.
The difference is: **+$170k**, 33.3% Higher.
Estimation is <u>wrongly outside 10%</u> of the error margin.

Your AI System prediction **with AI Explainer**: **$1,110k**.
The difference is: **+$599k**, 117% Higher.
Estimation is <u>wrongly outside 10%</u> of the error margin.

**Figure A.18: Review of performance on forward simulation trials, for participants to strengthen their understanding. The displayed UI is the Incremental condition.**

## [TASK] Describe the General AI System Behaviour

For all houses <u>overall</u>, what are the factors for each attribute? Please take note if it is positive or negative.

**Hint**:

1. Although you are not shown any specific houses, you can write <u>custom attribute values</u> by clicking on the values boxes.
2. You can input <u>custom factors</u> and the interface will compute the AI Explainer output automatically.
3. Try your factors on various houses with different attributes and observe if the outputs look accurate for <u>all houses overall</u>.



### Bathrooms



| -300k | -240k | -180k | -120k | -60k | 0k | 60k | 120k | 180k | 240k | 300k |

### Living Area



| -300k | -240k | -180k | -120k | -60k | 0k | 60k | 120k | 180k | 240k | 300k |

### Grade



| -300k | -240k | -180k | -120k | -60k | 0k | 60k | 120k | 180k | 240k | 300k |

### Age



| -300k | -240k | -180k | -120k | -60k | 0k | 60k | 120k | 180k | 240k | 300k |

### Adjustment



| -2,000k | -1,600k | -1,200k | -800k | -400k | 0k | 400k | 800k | 1,200k | 1,600k | 2,000k |

Do you agree or disagree with the following statement?

| | Strongly disagree | Disagree | Somewhat disagree | Neither agree nor disagree | Somewhat agree | Agree | Strongly agree |
|---|---|---|---|---|---|---|---|
| I am <u>confident</u> that the factors estimated above are correct. | ○ | ○ | ○ | ○ | ○ | ○ | ○ |

**Figure A.19: Explanation recall task for all instances overall (Global).**

# [TASK] Describe the AI System Behaviour for Living Area < 2.5

For the special groups of houses when Living Area < 2.5 ksqft, what are the factors for all attributes? Please take note if it is positive or negative.

**Hint**:

1. Although you are not shown any specific houses, you can write custom attribute values by clicking on the values boxes.
2. You can input custom factors and the interface will compute the AI Explainer output automatically.
3. Try your factors on various houses with different attributes and observe if the outputs look accurate for all houses overall.

| Attributes | Values | | Factors when Living Area < 2.5 | | Partial Price ($) |
|---|---|---|---|---|---|
| # Bathrooms | 0 | × | 0 | k = | ? |
| Living Area (ksqft) | 0 | × | 0 | k = | ? |
| Grade | 0 | × | 0 | k = | ? |
| Age (years) | 0 | × | 0 | k = | ? |
| Adjustment | | + | 0 | k = | ? |

Reset

AI Explainer ($): ? ≈ AI System ($): ?

**Bathrooms**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Living Area**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Grade**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Age**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Adjustment**

-2000  -1600  -1200  -800  -400  0  400  800  1200  1600  2000

Do you agree or disagree with the following statement?

| | Strongly disagree | Disagree | Somewhat disagree | Neither agree nor disagree | Somewhat agree | Agree | Strongly agree |
|---|---|---|---|---|---|---|---|
| I am confident that the factors estimated above are correct. | ○ | ○ | ○ | ○ | ○ | ○ | ○ |

**Figure A.20: Explanation recall task for subspace Living Area < 2.5.**

## [TASK] Describe the AI System Behaviour for Living Area ≥ 2.5

For the special groups of houses when Living Area ≥ 2.5 ksqft, what are the factors for all attributes? Please take note if it is positive or negative.

**Hint**:

1. Although you are not shown any specific houses, you can write custom attribute values by clicking on the values boxes.
2. You can input custom factors and the interface will compute the AI Explainer output automatically.
3. Try your factors on various houses with different attributes and observe if the outputs look accurate for all houses overall.

| Attributes | Values | Factors when Living Area ≥ 2.5 | Partial Price ($) |
|---|---|---|---|
| # Bathrooms | 0 | × 0 k = | ? |
| Living Area (ksqft) | 2.5 | × 0 k = | ? |
| Grade | 0 | × 0 k = | ? |
| Age (years) | 0 | × 0 k = | ? |
| Adjustment | | + 0 k = | ? |

Reset

AI Explainer ($) ? Infinity% Lower

AI System ($) 0k

**Bathrooms**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Living Area**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Grade**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Age**

-300  -240  -180  -120  -60  0  60  120  180  240  300

**Adjustment**

-2000  -1600  -1200  -800  -400  0  400  800  1200  1600  2000

Do you agree or disagree with the following statement?

| | Strongly disagree | Disagree | Somewhat disagree | Neither agree nor disagree | Somewhat agree | Agree | Strongly agree |
|---|---|---|---|---|---|---|---|
| I am confident that the factors estimated above are correct. | ○ | ○ | ○ | ○ | ○ | ○ | ○ |

**Figure A.21: Explanation recall task for subspace Living Area ≥ 2.5.**