



International Journal of Multidisciplinary Research and Growth Evaluation.

An Adaptive Explainability Framework for Machine Learning Predictions of Deals in Cloud Computing

Pavan Nithin Mullapudi

Senior Applied Scientist, Amazon Web Services, Seattle, WA, Independent Researcher, USA

* Corresponding Author: **Pavan Nithin Mullapudi**

Article Info

ISSN (online): 2582-7138

Volume: 05

Issue: 03

May-June 2024

Received: 17-05-2024

Accepted: 12-06-2024

Page No: 1027-1034

Abstract

Machine learning explainability frameworks often provide static, predefined explanations that fail to adapt to diverse user needs. This paper introduces a novel approach that enables users to personalize explanations through custom feature groupings, explanation goals, and context-specific preferences. Our framework allows users to logically organize model features into meaningful business-driven categories while defining the quality metrics that matter most to them. We introduce evaluation metrics including stability and diversity indices that align with user-defined objectives, bridging the gap between technical explainability and practical decision-making. Experiments with public datasets demonstrate that personalized explanations significantly improve decision-making confidence and user satisfaction compared to static approaches. The proposed framework transforms explanation from a static output to a dynamic, personalized process that adapts to specific user needs and contexts.

DOI: <https://doi.org/10.54660/IJMRGE.2024.5.3.1027-1034>

Keywords: Machine learning, predefined explanations, personalized process, technical explainability

1. Introduction

1.1 The Explainability Challenge

Explainable Artificial Intelligence (XAI) has emerged as a critical field for enhancing transparency and interpretability of machine learning models. As organizations across industries increasingly rely on complex ML models to drive decisions, the need for effective explanations becomes paramount. However, existing XAI frameworks typically provide explanations that are static and predetermined, failing to account for the diverse needs and preferences of different users.

The challenge is particularly acute in domains where various stakeholders interact with the same ML models but have fundamentally different information needs, technical expertise, and decision-making contexts. A data scientist may require detailed feature importance values, while a business analyst might need explanations framed in business terminology, and an executive may want high-level insights focused on strategic implications.

Traditional explainability approaches follow a one-size-fits-all methodology, where the same explanation method and presentation format is applied regardless of the user's role, knowledge level, or specific goals. This approach is fundamentally limited because it fails to recognize that explanation effectiveness is inherently contextual and subjective, depending heavily on who is consuming the explanation and why they need it.

1.2 The Need for Personalization

Several studies have focused on improving model interpretability through various XAI techniques. However, most of these approaches lack personalization capabilities related to how explanations are organized, presented, and evaluated. Users are treated as passive consumers of explanations rather than active participants in the explanation process.

This paper presents a novel framework addressing this gap by allowing users to define their explanation preferences based on specific needs and requirements. Our approach introduces three fundamental elements for enhancing personalization:

1. **Custom Feature Groupings:** Allowing users to logically group model features into meaningful categories that align with their mental models and domain understanding.
2. **Context Segmentation:** Enabling users to define specific contexts or segments where different explanation approaches might be appropriate.
3. **Explanation Goals:** Providing mechanisms for users to specify what makes an explanation valuable to them through customizable quality metrics.

By empowering users to personalize these elements, our framework transforms them from passive recipients to active participants in the explanation process. This user-centric approach ensures explanations are not only technically accurate but also relevant, meaningful, and actionable for each specific user.

1.3 Key Contributions

This paper makes the following key contributions:

1. We introduce a comprehensive framework for personalizing ML explanations through user-defined feature groupings, contexts, and quality metrics.
2. We propose novel evaluation metrics that capture explanation stability and diversity, enabling users to define explanation quality according to their specific needs.
3. We design and implement a user-centered interface that allows users to specify their explanation preferences intuitively.
4. We present experimental results demonstrating significant improvements in explanation utility and user satisfaction when explanations are personalized.
5. We provide practical guidelines for implementing personalized explanation systems across different domains and use cases.

The remainder of this paper is organized as follows: Section 2 reviews related work in explainable AI and user-centered explanation approaches. Section 3 details our methodology, including the system architecture and key components. Section 4 introduces our novel evaluation metrics for measuring explanation quality. Section 5 presents experimental results from applying our framework to real-world datasets. Section 6 discusses practical applications and implementation considerations. Section 7 addresses ethical considerations, and Section 8 concludes with a summary and directions for future research.

2. Related Work

2.1 Explainable AI Methods

Numerous methods have been suggested in recent literature to produce post hoc explanations for individual predictions made by complex ML models. Some of these local explanation techniques provide the influence of each feature on the model's prediction and are commonly referred to as local feature attribution methods.

Local feature attribution methods like LIME (Local Interpretable Model-agnostic Explanations) by Ribeiro et al. (2016) and SHAP (SHapley Additive exPlanations) by Lundberg and Lee (2017) have gained significant popularity due to their theoretical foundations and model-agnostic nature. LIME approximates a complex model locally with a simpler, interpretable model around a specific prediction.

SHAP assigns each feature an importance value representing its contribution to the prediction based on cooperative game theory.

Other approaches include counterfactual explanations (Wachter et al., 2017), which identify minimal changes to input features that would result in a different prediction, and example-based explanations (Kim et al., 2016), which explain predictions by referencing similar examples from the training data.

Their versatility has led to increasing utilization of feature attribution methods to explain complex models in various domains, including healthcare, finance, legal, and scientific applications. Consequently, ensuring the reliability of the explanations generated by these methods is critical to provide credible information about the underlying models to relevant stakeholders and decision-makers.

2.2 Evaluation of Explanations

Previous research has examined several aspects of explanation reliability. These include:

1. **Faithfulness (or fidelity):** This measures how accurately the explanation reflects the underlying model's behavior. Metrics include local fidelity (how well the explanation approximates the model locally) and global fidelity (how well the explanation captures overall model behavior).
2. **Stability (or robustness):** This assesses how consistent explanations are when inputs are slightly perturbed. Alvarez-Melis and Jaakkola (2018) visualized explanations generated by popular gradient-based methods and showed that they often lack robustness against small input perturbations.
3. **Fairness:** This evaluates whether explanations exhibit disparities across different demographic groups or data segments. Dai et al. (2022) proposed metrics to quantify disparities in explanation quality across groups.

Various evaluation metrics have been proposed to quantify these aspects, including probability changes (Slack et al., 2020), stability measurements (Alvarez-Melis and Jaakkola, 2018), and fairness metrics (Dai et al., 2022).

2.3 User-Centered Explanation Approaches

In a broader view, the aim of explanations is to make the reasons behind a decision or recommendation understandable to humans. As a result, the field of Explainable AI must adopt a human-centered approach. The Human-Computer Interaction (HCI) community has emphasized the importance of interdisciplinary collaboration and user-centered approaches in XAI (Abdul et al., 2018; Wang et al., 2019).

Currently, explanations are commonly designed based on the developer's intuition of what constitutes a "good explanation," often overlooking the perspective of the end-user. This approach neglects the diverse needs and preferences of different users, resulting in explanations that may be technically accurate but practically ineffective.

Research on user-centered explanations has demonstrated the benefits of explanations in various ways. These include providing the reasoning behind recommendations, enhancing user acceptance by presenting both positive and negative consequences, assisting users in making well-informed

decisions, and enabling effective communication between system providers and users.

Despite considerable research on generating and presenting explanations, the perspective of the end-user and their specific needs for explanations in different contexts remain relatively under-explored. Providing various types and levels of explanations could significantly impact the user's perception of the system, as the lack of appropriate explanations might lead to difficulty in understanding recommendations and negatively affect overall acceptance.

2.4 Gaps in Current Approaches

While significant progress has been made in developing various XAI methods and evaluation metrics, several important gaps remain:

1. **Static Explanation Methods:** Most current approaches use a single, predefined explanation method for all users and contexts, failing to adapt to different user needs.
2. **Limited User Input:** Users typically cannot specify how they want explanations organized or what metrics matter most to them.
3. **Technical Focus:** Explanations often emphasize technical details rather than aligning with users' mental models and domain knowledge.
4. **Rigid Feature Representation:** Features are typically presented individually without logical grouping into higher-level concepts that might be more meaningful to users.
5. **One-Dimensional Evaluation:** Explanation quality is often assessed using a single metric rather than considering multiple dimensions that might matter differently to different users.

Our framework addresses these gaps by providing a comprehensive approach to personalizing explanations through user-defined feature groupings, contexts, and quality metrics.

3. Methodology

3.1 System Architecture

Our framework for personalized explanations consists of four main components that work together to deliver customized explanations based on user preferences:

1. **User Preference Interface:** This component allows users to specify their preferences for explanations, including:
 - How they want features grouped into meaningful categories
 - What contexts or segments they are interested in
 - What quality metrics they prioritize for explanations
2. **Explanations Generator:** This component produces multiple explanations for each prediction using various XAI methods (e.g., LIME, SHAP, counterfactuals). Each method may provide different insights into the model's behavior, capturing different aspects of the prediction process.
3. **Explanation Metrics Calculator:** This component evaluates the quality of explanations across different dimensions, including:
 - **Stability:** How consistent explanations are across similar instances or slight variations in the model
 - **Diversity:** How much information content or variation exists in the explanations

- **Accuracy:** How faithfully the explanations reflect the model's actual behavior
- **Simplicity:** How easy the explanations are to understand

4. **Explanation Allocator:** This component matches the best explanation to each user based on their specified preferences and the quality metrics calculated for different explanation types.

The high-level design of the system is outlined in Figure 1. The framework integrates with existing ML pipelines, taking a trained model and predictions as inputs and generating personalized explanations as outputs.

[Figure 1: System Architecture diagram showing the flow from User Preferences through Explanation Generation, Metrics Calculation, and Allocation]

3.2 Custom Feature Groupings

A key innovation in our framework is allowing users to define custom groupings of model features that align with their mental models and domain expertise. Many ML models use dozens or hundreds of individual features, which can overwhelm users when presented in explanations. By enabling logical grouping of related features, our approach makes explanations more intuitive and actionable.

Users can create feature groupings in several ways:

1. **Semantic Groupings:** Features related to the same concept or entity can be grouped together. For example, in a customer churn prediction model, features related to usage patterns, support interactions, and billing history might form separate logical groups.
2. **Business Process Groupings:** Features can be organized according to their role in business processes. For instance, in a loan approval model, features might be grouped into application details, credit history, financial stability, and existing obligations.
3. **Actionability Groupings:** Features can be grouped based on whether and how they can be influenced or acted upon. Some features might be immutable (e.g., historical data), while others might be actionable through specific interventions.
4. **Hierarchical Groupings:** Features can be organized into hierarchies, allowing users to examine explanations at different levels of granularity. This enables both high-level overview explanations and detailed drill-down when needed.

The system stores these user-defined groupings and applies them when generating and presenting explanations. This approach transforms the explanation from a technical feature listing to a conceptually organized view that aligns with the user's understanding of the domain.

3.3 Context Segmentation

Different explanation approaches may be appropriate in different contexts or for different segments of data. Our framework allows users to define contexts where they want specific types of explanations, enabling context-sensitive personalization.

Context segmentation can be based on various factors:

1. **Data Segments:** Users can define segments based on data characteristics, such as customer segments, product categories, or time periods.

2. **Prediction Properties:** Contexts can be defined based on properties of the prediction itself, such as prediction confidence, prediction value range, or model uncertainty.
3. **Business Scenarios:** Users can specify different business scenarios or use cases where they need different explanation approaches.
4. **User Roles:** Different explanation contexts can be defined for different user roles, such as data scientists, business analysts, or executives.

By allowing context-specific explanation preferences, our framework ensures that explanations are not only personalized to the user but also adapted to the specific situation where they are needed.

3.4 Explanation Goals and Quality Metrics

Users have different goals when seeking explanations, and these goals should determine how explanation quality is defined and measured. Our framework allows users to specify their explanation goals and maps these to appropriate quality metrics.

Common explanation goals include:

1. **Decision Support:** Understanding why a specific prediction was made to inform a decision.
2. **Model Improvement:** Identifying potential issues or biases in the model to guide refinement.
3. **Process Compliance:** Verifying that the model's decision process complies with regulations or business rules.
4. **Knowledge Discovery:** Gaining new insights about patterns or relationships in the data.
5. **Trust Building:** Increasing confidence in the model's reliability and fairness.

Each of these goals may prioritize different aspects of explanation quality. For example, decision support may emphasize actionability and simplicity, while model improvement might prioritize faithfulness and comprehensiveness.

Our framework maps these high-level goals to specific quality metrics that can be measured and optimized, as detailed in the next section.

4. Explanation Quality Metrics

To enable personalization based on user goals, we define several metrics that capture different aspects of explanation quality. These metrics provide a quantitative basis for selecting the most appropriate explanation for each user and context.

4.1 Stability Metric

Stability measures how consistent explanations are across similar instances or slight variations in the model. This is particularly important for users who need reliable explanations that don't change dramatically with minor input changes.

We define stability through the concepts of local stability and global stability:

Local Stability Calculation:

1. **Base Explanation Calculation:** We calculate an explanation for a prediction using the complete training dataset, capturing the direction and magnitude of each feature or feature group's contribution.
2. **Modified Explanations Calculation:** We create multiple variations by retraining the model on different subsets of the training data or with slight perturbations to the input, generating multiple explanation variations.
3. **Local Stability Measurement:** The local stability of a feature or feature group is determined by measuring the proportion of variations where its contribution remains consistent in direction and similar in magnitude.

Global Stability:

Global stability is calculated as the average local stability across all features or feature groups. It provides an overall measure of how stable the explanations are at a broader level. A high stability score indicates that the explanations are robust to small changes in the training data or input, providing users with confidence that the explanation reflects consistent patterns rather than artifacts of specific data points.

4.2 Diversity Index

While stability is important, a model that always provides the same explanation regardless of the input would have high stability but low informational value. To capture the richness and usefulness of explanations, we introduce a diversity index.

The diversity index measures the variation and information content in explanations across a defined context or segment. It quantifies the extent to which different combinations of feature contributions are represented in the explanations.

We calculate the diversity index using Shannon's entropy concept, which measures the information content in a system: Diversity Index = Entropy / Max Possible Entropy where:

- Entropy is calculated as $-\sum P(x_i) * \log_2(P(x_i))$, with $P(x_i)$ being the probability of observing explanation pattern x_i
- Max Possible Entropy is the theoretical maximum entropy if all possible explanation patterns were equally likely

A higher diversity index indicates that the explanations provide varied insights across different instances, making them more informative for understanding the model's behavior in different situations.

4.3 Simplicity vs. Comprehensiveness

Different users may prefer explanations with different levels of detail. Some need comprehensive explanations that cover all relevant factors, while others prefer simpler explanations focused on the most important factors.

We measure simplicity using

1. **Feature Count:** The number of features or feature groups included in the explanation. Fewer features generally indicate higher simplicity.
2. **Concentration Ratio:** The proportion of the total impact captured by the top N features. A higher concentration indicates that a small number of features account for

most of the prediction.

3. **Cognitive Complexity:** A measure of how complex the relationships are among the features in the explanation, based on interaction effects and non-linear relationships.

Users can specify their preference along the simplicity-comprehensiveness spectrum, and the system will select explanations that align with this preference.

4.4 Actionability

For many users, the ultimate goal of an explanation is to inform actions. Actionability measures how well an explanation highlights factors that can be influenced or acted upon.

We assess actionability by:

1. **Mutable Feature Emphasis:** The proportion of the explanation focused on features that can be changed or influenced.
2. **Intervention Clarity:** How clearly the explanation indicates what changes would lead to different outcomes.
3. **Counterfactual Relevance:** Whether counterfactual explanations highlight realistic and meaningful alternative scenarios.

Users whose primary goal is to determine actions based on the model's predictions would prioritize explanations with high actionability scores.

5. Experimental Results

5.1 Experimental Setup

To evaluate our personalized explanation framework, we conducted experiments using several public datasets and ML models. For this paper, we focus on results from:

1. A customer churn prediction dataset with 21 features
2. A loan approval dataset with 17 features
3. A product recommendation dataset with 32 features

For each dataset, we trained gradient boosting models (specifically XGBoost) as these are commonly used in production environments and provide good prediction performance while being compatible with most explanation methods.

We implemented multiple explanation methods, including LIME, SHAP (both KernelSHAP and TreeSHAP variants), and counterfactual explanations. For each prediction, all methods were used to generate explanations, which were then evaluated according to the quality metrics described in Section 4.

We recruited 45 participants with varying roles (data scientists, business analysts, and domain experts) to use the system. Participants defined their custom feature groupings, specified contexts of interest, and indicated their explanation goals. They then evaluated the personalized explanations provided by our system compared to standard, non-personalized explanations.

5.2 Comparison of Explanation Methods

We first compared different explanation methods using our quality metrics to understand their strengths and weaknesses. Table 1 presents the results for global stability and diversity index across different methods.

[Table 1: Comparison of explanation methods across quality metrics]

For all datasets, SHAP methods generally provided higher stability than LIME, with TreeSHAP showing the highest stability scores. This indicates that SHAP explanations tend to be more consistent across similar instances or slight variations in the model.

However, LIME often produced higher diversity scores, particularly for the customer churn and product recommendation datasets. This suggests that LIME explanations capture more varied patterns across different instances, potentially providing richer insights into the model's behavior in different situations.

Counterfactual explanations showed moderate stability but high actionability scores, making them particularly suitable for decision support contexts where users need to understand what changes would lead to different outcomes.

5.3 Impact of Feature Groupings

We next evaluated how custom feature groupings affected user satisfaction and decision quality. Participants were asked to complete decision tasks using both standard feature-level explanations and explanations with their custom groupings.

Results showed that

1. **Comprehension Speed:** Users understood explanations with custom groupings 37% faster than explanations with individual features.
2. **Decision Confidence:** Users reported 42% higher confidence in decisions made using custom-grouped explanations.
3. **Decision Quality:** When evaluated against objective metrics, decisions made using custom-grouped explanations were 18% more aligned with optimal outcomes.

These results demonstrate that organizing features into meaningful groups significantly improves explanation utility. The benefits were particularly pronounced for business users and domain experts, who typically think in terms of business concepts rather than individual data features.

5.4 Context-Specific Explanations

We also examined how different explanation methods performed across user-defined contexts. Figure 2 shows the performance of different methods across contexts defined by prediction confidence levels.

[Figure 2: Performance of explanation methods across confidence contexts]

The results revealed interesting patterns:

1. For high-confidence predictions, stability was consistently high across methods, but SHAP provided slightly better performance.
2. For borderline predictions (those near the decision threshold), diversity and actionability became more important, with counterfactual methods showing advantages.
3. For low-confidence predictions, users valued explanations with higher diversity and comprehensiveness, as these helped understand the uncertainty in the prediction.

These findings highlight the importance of context-specific explanation selection. No single method performs best across all contexts, underscoring the value of our personalized approach.

5.5 User Satisfaction and Trust

Finally, we measured how personalized explanations affected user satisfaction and trust in the ML system. Participants reported:

1. **Overall Satisfaction:** 76% higher satisfaction with personalized explanations compared to standard explanations.
2. **System Trust:** 53% increase in trust in the ML system when using personalized explanations.
3. **Intention to Use:** 82% higher likelihood of using the ML system regularly when personalized explanations were available.

These results demonstrate that personalization significantly enhances the user experience and acceptance of ML systems. By providing explanations that align with users' mental models, preferences, and goals, our framework makes ML systems more accessible and trustworthy.

6. Practical Applications

The personalized explanation framework has broad applications across domains where ML is used for decision support. Here we highlight several practical applications and implementation considerations.

6.1 Customer Relationship Management

In CRM systems, ML models often predict customer behaviors such as churn, upsell potential, or lifetime value. Different stakeholders need different explanations:

- **Customer Service Representatives** might group features by actionable interventions they can offer, prioritizing simplicity and actionability.
- **Retention Specialists** might group features by customer pain points, valuing diversity to understand various churn drivers.
- **Product Managers** might group features by product usage patterns, emphasizing stability to identify consistent improvement opportunities.

By personalizing explanations for each role, organizations can ensure that insights from ML models translate into appropriate actions at every customer touchpoint.

6.2 Financial Services

In financial applications such as loan approval, investment recommendations, or fraud detection, explanations must balance technical accuracy with regulatory compliance:

- **Loan Officers** might group features by financial health indicators, prioritizing actionability to guide applicants on improving their chances.
- **Compliance Officers** might group features by regulatory categories, valuing stability to ensure consistent application of criteria.
- **Risk Analysts** might group features by risk factor types, emphasizing diversity to capture the full risk landscape.

Personalized explanations help financial institutions leverage ML effectively while maintaining regulatory compliance and customer trust.

6.3 Healthcare Decision Support

ML models in healthcare predict patient outcomes, recommend treatments, or identify high-risk individuals. Explanation needs vary widely:

- **Physicians** might group features by medical systems or conditions, valuing comprehensiveness and clinical relevance.
- **Patients** might prefer simpler groupings focused on lifestyle factors they can influence, prioritizing actionability.
- **Administrators** might group features by resource utilization categories, emphasizing diversity to understand various cost drivers.

By tailoring explanations to these different stakeholders, healthcare organizations can improve clinical decision-making, patient engagement, and operational efficiency.

6.4 Implementation Considerations

Organizations implementing personalized explanation systems should consider:

1. **User Research:** Conduct thorough research to understand how different users conceptualize domain concepts and what explanation goals they prioritize.
2. **Balanced Flexibility:** Provide suggested groupings and templates based on roles while allowing customization, balancing ease of use with personalization.
3. **Computing Constraints:** Generate and store multiple explanation types efficiently, possibly using asynchronous processing for computationally intensive methods.
4. **Feedback Mechanisms:** Implement ways for users to provide feedback on explanation quality, enabling continuous improvement of the system.
5. **Knowledge Sharing:** Create mechanisms for users to share effective feature groupings and explanation preferences, spreading best practices throughout the organization.

By addressing these considerations, organizations can maximize the value of personalized explanations while managing implementation complexity.

7. Ethical Considerations

The personalization of explanations raises important ethical considerations that must be addressed:

7.1 Balancing Personalization and Objectivity

While personalizing explanations improves user experience and understanding, it's essential to ensure that the underlying truths about model behavior aren't distorted. Personalization should change how information is organized and presented, not what information is conveyed.

Mitigation strategies include:

1. Distinguishing between factual aspects of explanations (feature contributions) and presentational aspects (grouping, emphasis)
2. Providing access to raw, unpersonalized explanations as a reference point
3. Monitoring for systematic biases introduced by personalization patterns

7.2 Preventing Confirmation Bias

There's a risk that users will create feature groupings that reinforce their existing beliefs rather than challenging them with new insights. This could lead to confirmation bias, where users only see explanations that align with their preconceptions.

To mitigate this risk:

1. Prompt users to consider alternative groupings or perspectives
2. Highlight unexpected or counterintuitive relationships in the data
3. Include diversity metrics that identify when explanations become too predictable

7.3 Transparency About Personalization

Users should understand that explanations are being personalized based on their preferences. Without this awareness, they might incorrectly assume that all users see the same explanations.

Best practices include:

1. Clearly indicating that explanations are personalized
2. Providing information about how personalization affects the explanations
3. Enabling users to view "standard" unpersonalized explanations for comparison

7.4 Equitable Access to Explanation Quality

There's a risk that personalization could create inequities in explanation quality, where some users receive more accurate or useful explanations than others based on their ability to configure the system effectively.

To ensure equitable access:

1. Provide high-quality default configurations for users who don't customize
2. Offer templates and presets based on role or use case
3. Monitor explanation quality across user groups to identify disparities

By addressing these ethical considerations proactively, organizations can implement personalized explanations in a responsible and beneficial manner.

8. Conclusion and Future Work**8.1 Summary of Contributions**

In this paper, we introduced a comprehensive framework for personalizing ML explanations through user-defined feature groupings, contexts, and quality metrics. Our approach

transforms users from passive recipients to active participants in the explanation process, ensuring explanations are not only technically accurate but also relevant, meaningful, and actionable.

Key contributions include:

1. A system architecture that integrates user preferences into the explanation generation process
2. Methods for users to define custom feature groupings that align with their mental models
3. Novel quality metrics including stability and diversity indices that capture different aspects of explanation quality
4. Empirical evidence demonstrating significant improvements in explanation utility and user satisfaction through personalization.

Our experiments across multiple datasets and user roles consistently showed that personalized explanations lead to faster comprehension, higher confidence, better decisions, and increased trust in ML systems.

8.2 Limitations

While our framework provides significant advances in explanation personalization, several limitations should be acknowledged:

1. **Computational Overhead:** Generating multiple explanation types for each prediction increases computational requirements, which may be challenging for real-time applications.
2. **User Effort:** Creating effective feature groupings requires initial effort from users, though this is mitigated by templates and sharing capabilities.
3. **Evaluation Complexity:** With personalized explanations, evaluation becomes more complex as different users may receive different explanations for the same prediction.
4. **Limited Model Coverage:** Our current implementation focuses on tabular data models; extending to other data types (text, images) would require additional research.

8.3 Future Research Directions

Our work opens several promising directions for future research:

1. **Automated Grouping Suggestions:** Developing methods to automatically suggest effective feature groupings based on data patterns and user behavior.
2. **Dynamic Adaptation:** Creating systems that learn user preferences over time and automatically adapt explanations based on usage patterns.
3. **Collaborative Explanations:** Exploring how multiple users with different perspectives can collaborate through shared explanations to gain comprehensive understanding.
4. **Multimodal Explanations:** Extending the framework to incorporate visual, textual, and interactive explanation elements based on user preferences.
5. **Longitudinal Studies:** Conducting long-term studies to understand how personalized explanations affect organizational decision-making and ML adoption over time.

By pursuing these research directions, we can continue to advance the state of the art in explanation personalization, making ML systems more transparent, trustworthy, and valuable to their users.

References

1. Abdul, A., Vermeulen, J., Wang, D., Lim, B. Y., & Kankanhalli, M. (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. Alvarez-Melis, D., & Jaakkola, T. S. (2018). On the Robustness of Interpretability Methods. arXiv:1806.08049.
3. Dai, J., Upadhyay, S., Aivodji, U., Bach, S. H., & Lakkaraju, H. (2022). Fairness via explanation quality: Evaluating disparities in the quality of post hoc explanations. In AAAI Conference on AI, Ethics, and Society.
4. Kim, B., Khanna, R., & Koyejo, O. O. (2016). Examples are not enough, learn to criticize! Criticism for interpretability. In Advances in Neural Information Processing Systems, 2280-2288.
5. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In Advances in Neural Information Processing Systems, 4765-4774.
6. Ribeiro, M. T., Singh, S., & Guestrin, C. (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.
7. Shannon, C. E. (1948). A Mathematical Theory of Communication. The Bell System Technical Journal, 27, 379-423.
8. Slack, D., Hilgard, S., Jia, E., Singh, S., & Lakkaraju, H. (2020). Fooling LIME and SHAP: Adversarial attacks on post hoc explanation methods. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, 180-186.
9. Wachter, S., Mittelstadt, B., & Russell, C. (2017). Counterfactual explanations without opening the black box: Automated decisions and the GDPR. Harvard Journal of Law & Technology, 31, 841-887.
10. Wang, D., Yang, Q., Abdul, A., & Lim, B. Y. (2019). Designing theory-driven user-centric explainable AI. In Proceedings of the 2019 CHI conference on human factors in computing systems, 1-15.