

SUFFICIENT AND NECESSARY EXPLANATIONS (AND WHAT LIES IN BETWEEN)

Anonymous authors

Paper under double-blind review

ABSTRACT

As complex machine learning models continue to be used in high-stakes decision settings, explaining their predictions is crucial. Post-hoc explanation methods aim to identify which features of an input \mathbf{x} are important to a model’s prediction $f(\mathbf{x})$. However, explanations often vary between methods and lack clarity, limiting the information we can draw from them. To address this, we formalize two precise concepts—*sufficiency* and *necessity*—to quantify how features contribute to a model’s prediction. We demonstrate that, although intuitive and simple, these two types of explanations may fail to fully reveal which features a model considers important. To overcome this, we propose and study a unified notion of importance that spans the entire necessity-sufficiency axis. Our unified notion, we show, has strong ties to other popular notions of feature importance, like those based on conditional independence and game-theoretic quantities like Shapley values. Lastly, through various experiments, we demonstrate that generating explanations along the necessity-sufficiency axis can uncover important features that may otherwise be missed and reveal that many post-hoc methods only provide features that are sufficient rather than necessary.

1 INTRODUCTION

Over recent years, modern machine learning (ML) models, mostly deep learning-based, have achieved impressive results across several complex domains. Models can now solve difficult image classification, inpainting, and segmentation problems, perform accurate text and sentiment analysis, predict the three-dimensional conformation of proteins, and more (LeCun et al., 2015; Wang et al., 2023). Despite their success, the rapid integration of these models into society requires caution (The White House, 2023). Modern ML systems are black-boxes, comprised of millions of parameters and non-linearities that obscure their prediction-making mechanisms from everyone. This lack of clarity raises concerns about explainability, transparency, and accountability (Zednik, 2021; Tomsett et al., 2018). Thus, understanding how these models work is essential for their safe deployment.

The lack of explainability has spurred research efforts in eXplainable AI (XAI), with a major focus on developing post-hoc methods to explain black-box model predictions, especially at a *local* level. For a model f and input $\mathbf{x} \in \mathbb{R}^d$, these methods aim to identify which features in \mathbf{x} are *important* for the model’s prediction, $f(\mathbf{x})$. They do so by estimating a notion of importance for each feature (or groups), which allows for a ranking of importance. Popular methods include CAM (Zhou et al., 2016), LIME (Ribeiro et al., 2016), gradient-based approaches (Selvaraju et al., 2017; Shrikumar et al., 2017; Jiang et al., 2021), rate-distortion techniques (Kolek et al., 2022), Shapley value-based explanations (Chen et al., 2018b; Teneggi et al., 2022; Mosca et al., 2022), perturbation-based methods (Fong & Vedaldi, 2017; Fong et al., 2019; Dabkowski & Gal, 2017), among others (Chen et al., 2018a; Yoon et al., 2018; Jethani et al., 2021; Wang et al., 2021; Ribeiro et al., 2018). However, many of these approaches lack rigor, as the meaning of their computed scores is often ambiguous. For example, it’s not always clear what large or negative gradients signify or what high Shapley values reveal about feature importance. To address these concerns, other research has focused on developing explanation methods based on logic-based definitions (Ignatiev et al., 2020; Darwiche & Hirth, 2020; Darwiche & Ji, 2022; Shih et al., 2018), conditional hypothesis testing (Teneggi et al. (2023); Tansey et al. (2022)), among formal notions. While these methods are a step towards rigor, they have drawbacks, including reliance on complex automated reasoners and limited ability to communicate their results in an understandable way for human decision-makers.

In this work, we advance XAI research by providing formal mathematical definitions of *sufficient* and *necessary* features for explaining complex ML models. First, we illustrate how, although informative, sufficient and necessary explanations offer incomplete insights into feature importance. To address this, we propose and study a more general unified framework for explaining models. Finally, we offer two novel perspectives on our framework through the lens of conditional independence and Shapley values, and crucially, show how it reveals new insights into feature importance.

1.1 SUMMARY OF OUR CONTRIBUTIONS

We propose and study two approaches, sufficiency, and necessity, which evaluate the contribution of a set of features in \mathbf{x} toward a model prediction $f(\mathbf{x})$. A sufficient set preserves the model’s output, while a necessary set, when removed, renders the output uninformative. Although the two concepts appear complementary, their precise relationship remains unclear. How similar are sufficient and necessary subsets? How different? To address these questions, we study the two concepts and propose a *unification* of both. Our contributions are summarized as follows:

1. We formalize precise mathematical definitions of sufficient and necessary features for model predictions that are related but complementary to those in previous works.
2. We propose a unified approach that combines sufficiency and necessity, exploring when and how they align or differ. Additionally, we motivate its utility by highlighting its connections to conditional independence and Shapley values, a game-theoretic measure of feature importance.
3. Through experiments of increasing complexity, we demonstrate how a unified perspective uncovers new, significant, and more comprehensive insights into feature importance.

2 SUFFICIENCY AND NECESSITY

Notation & Setting. We use boldface uppercase letters to denote random vectors (e.g., \mathbf{X}) and lowercase for their values (e.g., \mathbf{x}). For a subset $S \subseteq [d] := \{1, \dots, d\}$, we denote its cardinality by $|S|$ and its complement $S^c = [d] \setminus S$. Subscripts index features; e.g., \mathbf{x}_S represents \mathbf{x} restricted to the entries indexed by S . We consider a supervised learning setting with an unknown distribution \mathcal{D} over features $\mathcal{X} \subseteq \mathbb{R}^d$ and labels $\mathcal{Y} \subseteq \mathbb{R}$. We assume access to a model $f : \mathcal{X} \mapsto \mathcal{Y}$ that was trained on samples from \mathcal{D} . For an input $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, the goal is to identify the important features in \mathbf{x} for the prediction $f(\mathbf{x})$. To define importance, we will use the average restricted prediction, $f_S(\mathbf{x}) = \mathbb{E}_{\mathbf{X}_{S^c} \sim \mathcal{V}_{S^c}} [f(\mathbf{x}_S, \mathbf{X}_{S^c})]$, where \mathbf{x}_S is fixed and \mathbf{X}_{S^c} is a random vector drawn from an arbitrary reference distribution \mathcal{V}_{S^c} (which may or may not depend on S^c). For example, two common choices are the marginal $\mathcal{V}_{S^c} = p(\mathbf{X}_{S^c})$ and conditional distribution $\mathcal{V}_{S^c} = p(\mathbf{X}_{S^c} | \mathbf{x}_S)$. This strategy, popularized in (Lundberg & Lee, 2017; Lundberg et al., 2020), allows us to query f , which only takes inputs in \mathbb{R}^d , and analyze its behavior when sets of features are retained or removed.

Definitions. We now present our proposed definitions of sufficiency and necessity. At a high level, these definitions were formalized to align with the following guiding principles:

- P1. S is sufficient if it is enough to generate the original prediction, i.e. $f_S(\mathbf{x}) \approx f(\mathbf{x})$.
- P2. S is necessary if we cannot generate the original prediction without it, i.e. $f_{S^c}(\mathbf{x}) \not\approx f(\mathbf{x})$.
- P3. The set $S = [d]$ should be maximally sufficient and necessary for $f(\mathbf{x})$.

The principles P1 and P2 are natural and agree with the logical notions of sufficiency and necessity. Furthermore, because the full set of features provides all the information needed to make the prediction $f(\mathbf{x})$, it should thus be regarded as maximally sufficient and necessary (P3). With these principles laid out, we now formally define sufficiency and necessity.

Definition 2.1 (Sufficiency). *Let $\epsilon \geq 0$ and let $\rho : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$ be a metric on \mathbb{R} . A subset $S \subseteq [d]$ is ϵ -sufficient with respect to a distribution \mathcal{V} for f at \mathbf{x} if*

$$\Delta_{\mathcal{V}}^{\text{sup}}(S, f, \mathbf{x}) \triangleq \rho(f(\mathbf{x}), f_S(\mathbf{x})) \leq \epsilon. \quad (1)$$

Furthermore, S is ϵ -super sufficient if all supersets $\tilde{S} \supseteq S$ are ϵ -sufficient.

This notion of sufficiency is straightforward and aligns with P1. A subset S is ϵ -sufficient with respect to reference distribution \mathcal{V} if, with \mathbf{x}_S fixed, the average restricted prediction $f_S(\mathbf{x})$ is within ϵ from the original $f(\mathbf{x})$. Furthermore, S is ϵ -super sufficient if $\rho(f(\mathbf{x}), f_S(\mathbf{x})) \leq \epsilon$ and, $\forall \tilde{S} \supseteq S$,

$\rho(f(\mathbf{x}), f_{\tilde{S}}(\mathbf{x})) \leq \epsilon$. Namely, including more features in S keeps $f_S(\mathbf{x})$ ϵ close to $f(\mathbf{x})$. Note this definition aligns with P3, since the set $S = [d]$ is 0-sufficient (maximally sufficient). To find a small sufficient subset S of small cardinality $\tau > 0$, we can solve the following optimization problem:

$$\arg \min_{S \subseteq [d]} \Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x}) \quad \text{subject to } |S| \leq \tau \quad (\text{P}_{\text{suf}})$$

We will refer to this problem as the *sufficiency problem*, or (P_{suf}) . Using analogous ideas, we also define necessity and formulate an optimization problem to find small necessary subsets.

Definition 2.2 (Necessity). Let $\epsilon \geq 0$ and denote $\rho : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$ to be metric on \mathbb{R} . A subset $S \subseteq [d]$ is ϵ -necessary with respect to a distribution \mathcal{V} for f at \mathbf{x} if

$$\Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x}) \triangleq \rho(f_{S^c}(\mathbf{x}), f_{\emptyset}(\mathbf{x})) \leq \epsilon. \quad (2)$$

Furthermore, S is ϵ -super necessary if all supersets $\tilde{S} \supseteq S$ are ϵ -necessary.

Here, a subset S is ϵ -necessary if marginalizing out the features in S with respect to \mathcal{V}_S , results in an average restricted prediction $f_{S^c}(\mathbf{x})$ that is ϵ close to $f_{\emptyset}(\mathbf{x})$ – the average baseline prediction of f over $\mathcal{V}_{[d]}$. Furthermore, S is ϵ -super necessary if $\rho(f_S(\mathbf{x}), f(\mathbf{x})) \leq \epsilon$ and, $\forall \tilde{S} \supseteq S$, ϵ -necessary. Note, our definition of necessity differs from alternatives (Dhurandhar et al., 2018; Pawelczyk et al., 2020) which state that S is necessary if $\rho(f(\mathbf{x}), f_{S^c}(\mathbf{x})) \geq \Delta$ for some $\Delta > 0$. Our notion is more general in that it implies this condition. Intuitively, if $f_{\emptyset}(\mathbf{x})$ and $f(\mathbf{x})$ differ, and $f_{S^c}(\mathbf{x})$ is close to $f_{\emptyset}(\mathbf{x})$, then $f_{S^c}(\mathbf{x})$ and $f(\mathbf{x})$ will also differ. Furthermore, for $S = [d]$, we have $\Delta^{\text{nec}}_{\mathcal{V}}(S, f, \mathbf{x}) \triangleq \rho(f_{\emptyset}(\mathbf{x}), f_{\emptyset}(\mathbf{x})) = 0$, indicating that $S = [d]$ is 0-necessary (maximally necessary) as desired. A detailed comparison of our approach with classical definitions, along with its advantages, is provided in the Appendix. To identify a ϵ -necessary subset S of small cardinality $\tau > 0$, one can solve the following optimization problem, which we refer to as the *necessity problem* or (P_{nec}) .

$$\arg \min_{S \subseteq [d]} \Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x}) \quad \text{subject to } |S| \leq \tau \quad (\text{P}_{\text{nec}})$$

Having presented our definitions, we now discuss related works before presenting our main results.

3 RELATED WORK

Notions of sufficiency, necessity, their duality and connections with other feature attribution methods have been studied to varying degrees. We comment on the main related works in this section.

Sufficiency. The notion of sufficient features has gained significant attention in recent research. Shih et al. (2018) explore a symbolic approach to explain Bayesian network classifiers and introduce prime implicant explanations, which are minimal subsets S that make features in the complement irrelevant to the prediction $f(\mathbf{x})$. For models represented by a finite set of first-order logic (FOL) sentences, Ignatiev et al. (2020) refer to prime implicants as abductive explanations (AXp’s). For classifiers defined by propositional formulas and inputs with discrete features, Darwiche & Hirth (2020) refer to prime implicants as sufficient reasons and define a complete reason to be the disjunction of all sufficient reasons. They present efficient algorithms, leveraging Boolean circuits, to compute sufficient and complete reasons and demonstrate their use in identifying classifier dependence on protected features that should not inform decisions. For more complex models, Ribeiro et al. (2018) propose high-precision probabilistic explanations called anchors, which represent local, sufficient conditions. For \mathbf{x} positively classified by f , Wang et al. (2021) propose a greedy approach to solve (P_{suf}) , I Amoukou & Brunel (2022) extend this work to regression settings using tree-based models, and Fong & Vedaldi (2017) introduce the preservation method which relaxes S to $[0, 1]^d$.

Necessity. There has also been significant focus on identifying necessary features – those that, when altered, lead to a change in the prediction $f(\mathbf{x})$. For models expressible by FOL sentences, Ignatiev et al. (2019) define prime implicants as the minimal subsets that when changed, modify the prediction $f(\mathbf{x})$ and relate these to adversarial examples. For Boolean models predicting on samples \mathbf{x} with discrete features, Ignatiev et al. (2020) and (Darwiche & Hirth, 2020) refer to prime implicants as contrastive explanations (CXp’s) and necessary reasons, respectively. Beyond boolean functions, for \mathbf{x} positively classified by a classifier f , Fong et al. (2019) relax S to $[0, 1]^d$ and propose the deletion method to approximately solve (P_{nec}) .

Duality Between Sufficiency and Necessity. Dabkowski & Gal (2017) characterize the preservation and deletion methods as discovering the *smallest sufficient* and *destroying region* (SSR and SDR).

They propose combining the two but do not explore how solutions to this approach may differ from individual SSR and SDR solutions. Ignatiev et al. (2020) show that AXp’s and CXp’s are minimal hitting sets of another by using a hitting set duality result between minimal unsatisfiable and correction subsets. The result enables the identification of AXp’s from CXp’s and vice versa.

Sufficiency, Necessity, and General Feature Attribution Methods. Precise connections between sufficiency, necessity, and other popular feature attribution methods (such as Shapley values (Shapley, 1951; Chen et al., 2018b; Lundberg & Lee, 2017)) remains unclear. To our knowledge, Covert et al. (2021) provide the only work examining these approaches (Fong & Vedaldi, 2017; Fong et al., 2019; Dabkowski & Gal, 2017) in the context of general removal-based methods, i.e., methods that remove certain input features to evaluate different notions of importance. The work of Watson et al. (2021) is also relevant to our work, as it formalizes a connection between notions of sufficiency and Shapley values. With the specific payoff function defined as $v(S) = \mathbb{E}[f(\mathbf{x}_S, \mathbf{X}_{S^c})]$, they show how each summand in the Shapley value measures the sufficiency of feature i to a particular subset.

4 UNIFYING SUFFICIENCY AND NECESSITY

Given a model f and sample \mathbf{x} , we can identify a small set of important features S by solving either (P_{suf}) or (P_{nec}) . While both methods are popular (Kolek et al., 2022; Fong & Vedaldi, 2017; Bhalla et al., 2023; Yoon et al., 2018), identifying small sufficient or necessary subsets may not provide a complete picture of how f uses \mathbf{x} to make a prediction. To see why, consider the following scenario: for a fixed $\tau > 0$, let S^* be a ϵ -sufficient solution to (P_{suf}) , so that $|S^*| \leq \tau$ and $\Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x}) \leq \epsilon$. While S^* is ϵ -sufficient, it can also be true that $\Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x}) > \epsilon$ indicating S^* is **not** ϵ -necessary: indeed, this can simply happen when its complement, S^{c*} , contains important features. This scenario raises two questions: 1) How different are sufficient and necessary features? 2) How does varying the levels of sufficiency and necessity affect the optimal set of important features?

To answer these important questions (and avoid the scenario above) we propose studying a unification of (P_{suf}) and (P_{nec}) . Consider $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) = \alpha \cdot \Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x}) + (1 - \alpha) \cdot \Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x})$, a convex combination of $\Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x})$ and $\Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x})$, where $\alpha \in [0, 1]$ controls the extent to which S is sufficient vs. necessary. Our *unified problem*, (P_{uni}) , can be expressed as:

$$\arg \min_{S \subseteq [d]} \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) \quad \text{subject to } |S| \leq \tau \quad (P_{\text{uni}})$$

When α is 1 or 0, $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$ reduces to $\Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x})$ or $\Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x})$, respectively. In these extreme cases, S is only sufficient or necessary. In the remainder of this work we will theoretically analyze (P_{uni}) , characterize its solutions, and provide different interpretations of what properties the solutions have through the lens of conditional independence and game theory. In the experimental section, we will show that solutions to (P_{uni}) provide insights that neither (P_{suf}) nor (P_{nec}) offer.

4.1 SOLUTIONS TO THE UNIFIED PROBLEM

We begin with a simple lemma that demonstrates why (P_{uni}) enforces both sufficiency and necessity.

Lemma 4.1. *Let $\alpha \in (0, 1)$. For $\tau > 0$, denote S^* to be a solution to (P_{uni}) for which $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) = \epsilon$. Then, S^* is $\frac{\epsilon}{\alpha}$ -sufficient and $\frac{\epsilon}{1-\alpha}$ -necessary. Formally,*

$$0 \leq \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{\alpha} \quad \text{and} \quad 0 \leq \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{1-\alpha}. \quad (3)$$

The proof of this result, and all others, is included Appendix A.1. This result illustrates that solutions to (P_{uni}) satisfy varying definitions of sufficiency and necessity. Furthermore, as α increases from 0 to 1, the solution shifts from being highly necessary to highly sufficient. In the following results, we will show *when* and *how* solutions to (P_{uni}) are similar (and different) to those of (P_{suf}) and (P_{nec}) . To start, we present the following lemma, which will be useful in subsequent results.

Lemma 4.2. *For $0 \leq \epsilon < \frac{\rho(f(\mathbf{x}), f_{\theta}(\mathbf{x}))}{2}$, denote S_{suf}^* and S_{nec}^* to be ϵ -sufficient and ϵ -necessary sets. Then, if S_{suf}^* is ϵ -super sufficient or S_{nec}^* is ϵ -super necessary, we have $S_{\text{suf}}^* \cap S_{\text{nec}}^* \neq \emptyset$.*

This lemma demonstrates that, given ϵ -sufficient and necessary sets S_{suf}^* and S_{nec}^* , if either additionally satisfies the stronger notions of super sufficiency or necessity, they must share some features. This proves useful in characterizing a solution to (P_{uni}) , which we now do in the following theorem.

Theorem 4.1. Let $\tau_1, \tau_2 > 0$ and $0 \leq \epsilon < \frac{1}{2} \cdot \rho(f(\mathbf{x}), f_\emptyset(\mathbf{x}))$. Denote S_{suf}^* and S_{nec}^* to be ϵ -super sufficient and ϵ -super necessary solutions to (P_{suf}) and (P_{nec}) , respectively, such that $|S_{\text{suf}}^*| = \tau_1$ and $|S_{\text{nec}}^*| = \tau_2$. Then, there exists a set S^* such that

$$\Delta_Y^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) \leq \epsilon \quad \text{and} \quad \max(\tau_1, \tau_2) \leq |S^*| < \tau_1 + \tau_2. \quad (4)$$

Furthermore, if $S_{\text{suf}}^* \subseteq S_{\text{nec}}^*$ or $S_{\text{nec}}^* \subseteq S_{\text{suf}}^*$ then $S^* = S_{\text{nec}}^*$ or $S^* = S_{\text{suf}}^*$ respectively.

This result demonstrates that when there are ϵ -super sufficient and ϵ -super necessary solutions to (P_{suf}) and (P_{nec}) , then one can identify a set S^* with small Δ^{uni} . As an example, consider features that are ϵ -super sufficient, S_{suf}^* . If we have domain knowledge that $S_{\text{suf}}^* \subseteq S_{\text{nec}}^*$, and S_{nec}^* is ϵ -super necessary, then S_{nec}^* will have a small Δ^{uni} . Conversely, if we know that S_{suf}^* is ϵ -super necessary along with being a subset of ϵ -super sufficient set S_{suf}^* , then S_{suf}^* will have a small Δ^{uni} .

5 TWO PERSPECTIVES OF THE UNIFIED APPROACH

In the previous section, we characterized solutions to (P_{uni}) and their connections to those of (P_{suf}) and (P_{nec}) . To further motivate and the unified approach, we now offer two alternative perspectives of our framework through the lens of conditional independence and Shapley values.

5.1 A CONDITIONAL INDEPENDENCE PERSPECTIVE

Here we demonstrate how sufficiency, necessity, and their unification, can be understood as conditional independence relations between features \mathbf{X} and label Y .

Corollary 5.1. Suppose $\forall S \subseteq [d]$, $\mathcal{V}_S = p(\mathbf{X}_S | \mathbf{X}_{S^c} = \mathbf{x}_{S^c})$. Let $\alpha \in (0, 1)$, $\epsilon \geq 0$, and denote $\rho : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$ to be a metric. Furthermore, for $\tau > 0$ and $f(\mathbf{X}) = \mathbb{E}[Y | \mathbf{X}]$, let S^* be a solution to (P_{uni}) such that $\Delta_Y^{\text{uni}}(S, f, \mathbf{x}, \alpha) = \epsilon$. Then, S^* satisfies the follow conditional independencies,

$$\rho(\mathbb{E}[Y | \mathbf{x}], \mathbb{E}[Y | \mathbf{X}_{S^*} = \mathbf{x}_{S^*}]) \leq \frac{\epsilon}{\alpha} \quad \text{and} \quad \rho(\mathbb{E}[Y | \mathbf{X}_{S^c} = \mathbf{x}_{S^c}], \mathbb{E}[Y]) \leq \frac{\epsilon}{1 - \alpha}. \quad (5)$$

The assumption in this corollary is that, $\forall S \subseteq [d]$, $f_S(\mathbf{x})$ is evaluated using the conditional distribution $p(\mathbf{X}_{S^c} | \mathbf{X}_S = \mathbf{x}_S)$ as the reference distribution \mathcal{V}_S . Given the recent advancements in generative models (Song & Ermon, 2019; Ho et al., 2020; Song et al., 2021), this assumption is (approximately) reasonable in many practical settings, as we will demonstrate in our experiments. For this particular \mathcal{V}_S , the result shows that minimizing (P_{uni}) with model $f(\mathbf{X}) = \mathbb{E}[Y | \mathbf{X}]$ identifies an S^* that approximately satisfies two conditional independence properties. First, S^* is sufficient as conditioning on S^* leaves the complement S^{c*} with minimal additional information about Y . Second, S^* is necessary because when we solely rely on the complement S^{c*} , the information gained about Y is minimal and similar to $\mathbb{E}[Y = 1]$.

5.2 A SHAPLEY VALUE PERSPECTIVE

In the previous section, we detailed the conditional independence relations being optimized for when solving (P_{uni}) . We now present an arguably less intuitive result that shows that solving (P_{uni}) is equivalent to maximizing the lower bound of the Shapley value. Before presenting our result, we provide a brief background on this game-theoretic quantity.

Shapley Values. Shapley values use game theory to measure the importance of players in a game. Let the tuple $([n], v)$ represent a cooperative game with players $[n] = \{1, 2, \dots, n\}$ and denote a characteristic function $v(S) : \mathcal{P}([n]) \rightarrow \mathbb{R}$. Then, the Shapley value (Shapley, 1951) for player j in the game $([n], v)$ is $\phi_j^{\text{shap}}([n], v) = \sum_{S \subseteq [n] \setminus \{j\}} w_S \cdot [v(S \cup \{j\}) - v(S)]$ where $w_S = \frac{|S|!(n-|S|-1)!}{n!}$. In the context of XAI, Shapley values are widely used to measure local feature importance by treating input features as players in a game (Covert et al., 2020; Teneggi et al., 2022; Chen et al., 2018b; Lundberg & Lee, 2017). Given a sample \mathbf{x} and a model f , the importance of x_j to the prediction $f(\mathbf{x})$ is measured by computing ϕ_j^{shap} for a game $([d], v)$, where $v(S)$ quantifies how the features in S contribute to $f(\mathbf{x})$. Different choices of $v(S)$ can be found in (Lundberg & Lee, 2017; Sundararajan & Najmi, 2020; Watson et al., 2024). Although computing ϕ_j^{shap} is computationally intractable, several practical methods for estimation have been developed (Chen et al., 2023; Teneggi et al., 2022; Zhang et al., 2023; Lundberg et al., 2020). While Shapley values are popular across various domains (Moncada-Torres et al., 2021; Zoabi et al., 2021; Liu et al., 2021), few works, aside from Watson et al. (2021), explore their connections to sufficiency and necessity.

With this background, we now present our result. Recall solving (P_{uni}) finds a small subset S with low $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$. Notice that (P_{uni}) naturally *partitions* the features into two sets, S and S^c . In the following theorem we demonstrate that finding a small S with minimal $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$ is equivalent to maximizing a lower bound on the Shapley value in a two player game.

Theorem 5.1. *Consider an input \mathbf{x} for which $f(\mathbf{x}) \neq f_{\emptyset}(\mathbf{x})$. Denote by $\Lambda_d = \{S, S^c\}$ the partition of $[d] = \{1, 2, \dots, d\}$, and define the characteristic function to be $v(S) = -\rho(f(\mathbf{x}), f_S(\mathbf{x}))$. Then,*

$$\phi_S^{\text{shap}}(\Lambda_d, v) \geq \rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x})) - \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha). \quad (6)$$

This result motivates minimizing $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$ via a game-theoretic interpretation. The tuple (Λ_d, v) specifies a game, and since there are 2^{d-1} ways to partition $[d]$ into 2 subsets, there are 2^{d-1} games. The inequality above holds for each of them. Thus, Theorem 5.1 implies that finding the S with minimal $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$ is equivalent to identifying the game (i.e. partition) (Λ_d, v) in which S has the largest lower bound on its Shapley value.

6 SOLVING THE UNIFIED PROBLEM

Before presenting our results, we briefly discuss different approaches to solving (P_{uni}) . In general, this problem is NP-hard however, in certain settings, one can efficiently compute exact solutions or use tractable relaxations, (Kolek et al., 2022; Fong et al., 2019; Linder et al., 2022) to approximate solutions. We present these general approaches here, and defer details to Appendix A.2.

Exhaustive Search. When the feature space dimension, d , or choice of $\tau \in \mathbb{Z}_{>0}$ is small an exhaustive search can compute exact solutions to (P_{uni}) by evaluating $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha)$ for all $\binom{d}{\tau}$ subsets S of cardinality τ and selecting the minimizer.

Instance-wise Optimization. When d is large, rendering (P_{uni}) intractable, one can generate approximate solutions by solving the relaxed problem¹

$$\arg \min_{S \subseteq [0,1]^d} \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) + \lambda_1 \cdot \|S\|_1 + \lambda_{\text{TV}} \cdot \|S\|_{\text{TV}}. \quad (7)$$

This type of approach is often used in computer vision and natural language problems (Fong et al., 2019; Kolek et al., 2022; Linder et al., 2022) to generate instance-specific solutions.

Parametric Model Approach. Another we approach we take to generate solutions to (P_{uni}) is to learn models $g_{\theta} : \mathcal{X} \mapsto [0, 1]^d$ that (approximately) solve the following optimization problem:

$$\arg \min_{\theta \in \Theta} \mathbb{E}_{\mathbf{X} \sim \mathcal{D}_{\mathcal{X}}} \left[\Delta_{\mathcal{V}}^{\text{uni}}(g_{\theta}(\mathbf{X}), f, \mathbf{X}, \alpha) + \lambda_1 \cdot \|g_{\theta}(\mathbf{X})\|_1 + \lambda_{\text{TV}} \cdot \|g_{\theta}(\mathbf{X})\|_{\text{TV}} \right]. \quad (8)$$

With these models, an approximate solution can be computed via $g_{\theta}(\mathbf{x})$. This method is popular (Chen et al., 2018a; Yoon et al., 2018; Linder et al., 2022), as it handles highly structured data well and requires training only one model, rather than repeatedly solving Eq. (7) for each sample.

7 EXPERIMENTS

We demonstrate our theoretical findings in multiple settings of increasingly complexity: two tabular data tasks (on synthetic data and the US adult income dataset (Ding et al., 2021)) and two high-dimensional image classification tasks using the RSNA 2019 Brain CT Hemorrhage Challenge (Flanders et al., 2020) and CelebA-HQ datasets (Lee et al., 2020)

7.1 TABULAR DATA

With the following tabular data settings, we demonstrate how the specific trade-off between sufficiency and necessity can greatly alter the solutions to (P_{uni}) . To do so, we compute exact solutions via exhaustive search to (P_{uni}) for varying levels of sufficiency vs. necessity and multiple size constraints. We learn a predictor f and, for 100 new samples, solve (P_{uni}) for $\tau \in \{3, 6, 9\}$ and $\alpha \in [0, 1]$, with $\rho(a, b) = |a - b|$ and $\mathcal{V}_S = p(\mathbf{X}_S | \mathbf{X}_{S^c} = \mathbf{x}_{S^c})$. For a fixed τ and sample \mathbf{x} , we denote $S_{\alpha_i}^*$ to be a solution to (P_{uni}) for α_i . It is represented as a binary vector $s \in \{0, 1\}^{10}$, where $s_j = 1$ if $j \in S_{\alpha_i}^*$ and 0 otherwise. To analyze the stability of $S_{\alpha_i}^*$ as sufficiency and necessity vary, we report the normalized average Hamming distance (Hamming, 1950) between $S_{\alpha_i}^*$ and S_0^* (with 95% confidence intervals) as a function of α .

¹Here, λ_1 , $\|S\|_1$ and λ_{TV} , $\|S\|_{\text{TV}}$ are the ℓ_1 and Total Variation norms and hyperparameters, respectively, promoting sparsity and smoothness.

7.1.1 LINEAR REGRESSION

We begin with a regression example. Features are distributed as $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}\mathbf{A}^T)$ with $\boldsymbol{\mu} = [2^i]_{i=1}^d$ and $\mathbf{A}_{i,j} \sim U(0, 1)$. The response is $Y = \boldsymbol{\beta}^T \mathbf{X} + \epsilon$, with $\boldsymbol{\beta} = 32 \cdot [2^{-i}]_{i=1}^d$ and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d \times d})$. We fix $d = 10$ and use the model $f(\mathbf{X}) = \hat{\boldsymbol{\beta}}^T \mathbf{X}$, where $\hat{\boldsymbol{\beta}}$ is the least squares solution.

Stability of Unified Solutions. Fig. 1a shows that when solutions are constrained to be small ($\tau = 3$), increasing α to enforce greater sufficiency results in a steady increase in Hamming distance, indicating that the solutions $S_{\alpha_i}^*$ are consistently changing. When larger solutions are allowed ($\tau = 6$), $S_{\alpha_i}^*$ rapidly changes with the introduction of sufficiency, as seen by the initial steep rise in Hamming distance. However, as α continues to increase, this distance grows more gradually. Lastly, when the solution size approaches the dimension of the feature space ($\tau = 9$), small to medium levels of sufficiency do not significantly alter $S_{\alpha_i}^*$. However, high levels of sufficiency ($\alpha > 0.8$) lead to extreme changes in the solutions, as shown by a sharp increase in Hamming distance.

7.1.2 AMERICAN COMMUNITY SURVEY INCOME (ACSIIncome)

We use the ACSIncome dataset for California, including 10 demographic and socioeconomic features such as age, education, occupation, and geographic region. We train a Random Forest classifier to predict whether an individual’s annual income exceeds \$50K, achieving a test accuracy $\approx 81\%$.

Stability of Unified Solutions. Fig. 1b shows that when solutions are forced to be small ($\tau = 3$), increasing α to enforce sufficiency results in a steady increase in Hamming distance, indicating the solutions $S_{\alpha_i}^*$ are changing. For larger solutions ($\tau = 6$), $S_{\alpha_i}^*$ changes significantly when low levels sufficiency are required, indicated by initial rise in the Hamming distance. As α continues to increase, the Hamming distance grows more gradually. Interestingly, when the size is close to feature space’s dimensionality ($\tau = 9$), the Hamming distance exhibits a behavior similar to that observed for $\tau = 3$. In conclusion, both examples show that the optimal feature set can vary depending on the size constraint and balance between sufficiency and necessity.

7.2 IMAGE CLASSIFICATION

The following two experiments explore high dimensional image classification tasks. The features are pixel values and so a subset S corresponds to a binary mask identifying important pixels. Since solving (P_{suf}) , (P_{nec}) , or (P_{uni}) is intractable here, we use two methods, the per-sample and model based approach in Eqs. (7) and (8) to identify sufficient and necessary masks. These experiments serve two purposes. First, they will analyze the ability popular explanation methods—including Integrated Gradients (Sundararajan et al., 2017), GradientSHAP (Lundberg & Lee, 2017), Guided GradCAM (Selvaraju et al., 2017), and h-Shap (Teneggi et al., 2022)—to identify small sufficient and necessary subsets. To ensure consistent analysis, all attribution scores are normalized to the interval $[0, 1]$. This is done by setting the top 1% of nonzero scores to 1 and dividing the remaining by the minimum score from the top 1% nonzero scores, which is common practice (Kokhlikyan et al., 2020). Binary masks are then generated by thresholding the normalized scores using thresholds $t \in (0, 1)$. For a test set of images and normalized attribution scores, we report the average (across all binary masks) $-\log(\Delta^{\text{suf}})$, $-\log(\Delta^{\text{nec}})$, and $-\log(L^0)$ where L^0 is the relative size of S for $t \in (0, 1)$ to analyze the sufficiency, necessity and size of the explanations. The second objective of these experiments is to understand and visualize the similarities and differences between sufficient and necessary sets.

7.2.1 RSNA CT HEMORRHAGE

We use the RSNA 2019 Brain CT Hemorrhage Challenge dataset comprised of 752,803 scans. Each scan is annotated by expert neuroradiologists with the presence and type(s) of hemorrhage (i.e., epidural, intraparenchymal, intraventricular, subarachnoid, or subdural). We use a ResNet18 (He et al., 2016) classifier that was pretrained on this data (Teneggi et al., 2022). Since the dataset

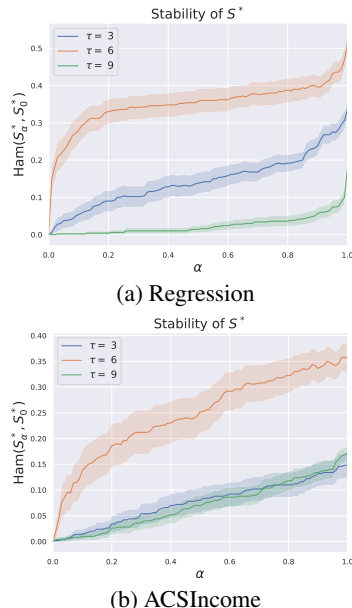


Figure 1: Stability of (P_{uni}) Solutions

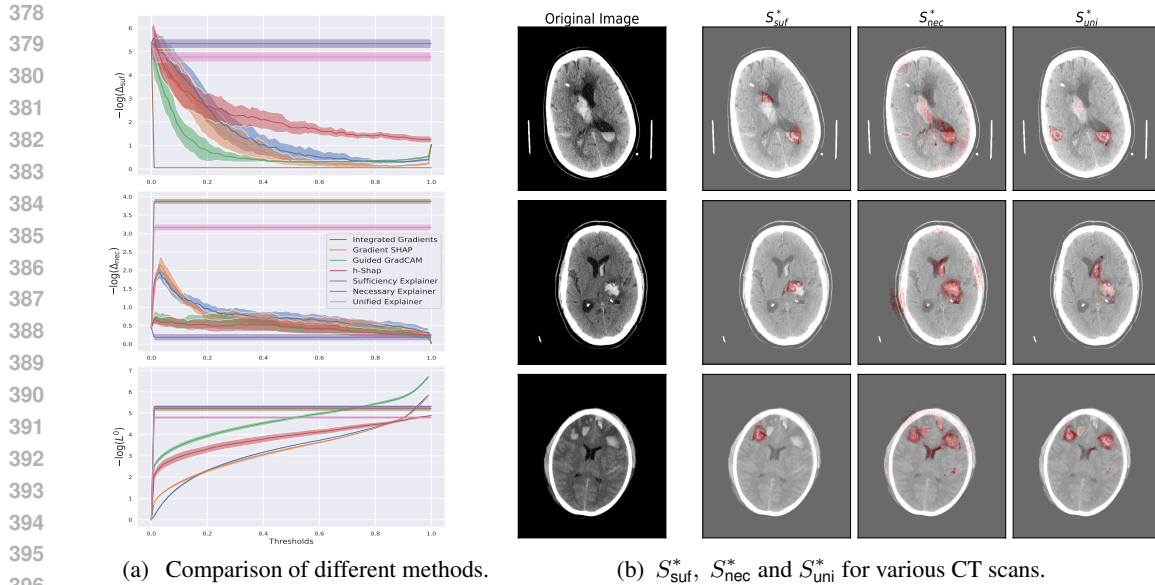


Figure 2: Experimental results on the RSNA dataset.

consists of highly complex and diverse images, we employ the per-example approach in Eq. (7) with $\alpha \in \{0, 0.5, 1\}$ to learn sufficient and necessary masks. Further details are in Appendix A.2.

Comparison of Post-hoc Interpretability Methods. For a set of 20 images positively classified by the ResNet model, we apply multiple post-hoc interpretability methods, as well as compute sufficient and necessary masks by solving (7). The results in Fig. 2a show that for thresholds $t < 0.1$, many methods identify sufficient sets smaller in size than the sufficient and unified explainer, as indicated by their large values of $-\log(\Delta^{suf})$ and smaller values of $-\log(L^0)$. However, for $t > 0.1$, only the sufficient and unified explainer identify sufficient sets of a constant small size. Importantly, *no methods, besides the necessity and unified explainers, identify necessary sets*. Furthermore, as expected, the sufficient explainer does not identify necessary sets and vice versa. The unified explainer, as expected, identifies a sufficient and necessary set (at the cost of a larger set). In conclusion, while off-the-shelf methods can identify sufficient, they do not identify necessary sets for small thresholds.

Sufficiency vs. Necessity. In Fig. 2b we visualize the sufficient and necessary features in various CT scans. The first observation is that sufficient subsets do not provide a complete picture of which features are important. Notice for all the CT scans, a sufficient set, S_{suf}^* highlights one or two, but never all, brain hemorrhages in the scans. For example, in the last row, S_{suf}^* only contains the right frontal lobe parenchymal hemorrhages, which happens to be one of the larger hemorrhages present. On the other hand, necessary sets, S_{nec}^* , contain parts of, sometimes entirely, *all* hemorrhages in the scans. In the last row, S_{nec}^* contains all multifocal parenchymal hemorrhages in both right and left frontal lobes, because when all these regions are masked, the model yields a prediction ≈ 0.64 —the prediction of the model on the mean image. Finally, notice in the 2nd and 3rd columns that S_{nec}^* and S_{uni}^* are nearly identical, which precisely demonstrate Lemma 4.1 and Theorem 4.1 in practice. First, since S_{suf}^* is super sufficient, S_{suf}^* and S_{nec}^* share common features. Second, visually $S_{suf}^* \subseteq S_{nec}^*$ holds approximately and so $S_{nec}^* = S_{uni}^*$. Through this experiment we are able to highlight the differences between sufficient and necessary sets, show how each contain important and complementary information, and demonstrate our theory holding in real world settings.

7.2.2 CELEBA-HQ

We use a modified version of the CelebA-HQ dataset (Karras, 2017) that contains 30,000 celebrity faces resized to 256×256 pixels. We train a ResNet18 to classify whether a celebrity is smiling, achieving a test accuracy $\approx 94\%$ and use the model based approach via solving Eq. (8) to generate sufficient and necessary masks. Given the structured nature of the dataset and the similarity of features across images, we use the model approach because it prevents overfitting to spurious signals

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

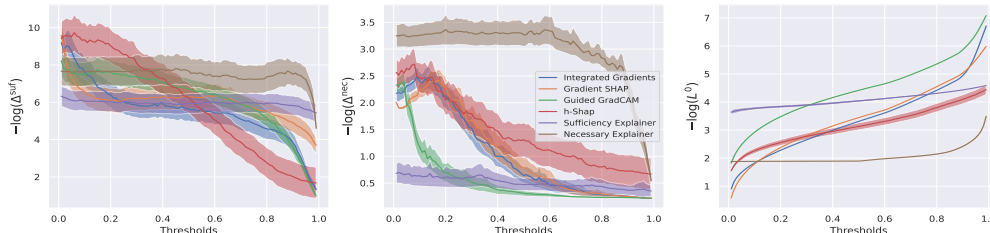


Figure 3: Comparison of different methods on the CelebAHQ dataset.

(Linder et al., 2022), an issue that can arise with per-example methods. Implementation details and hyperparameter settings are included in Appendix A.2.

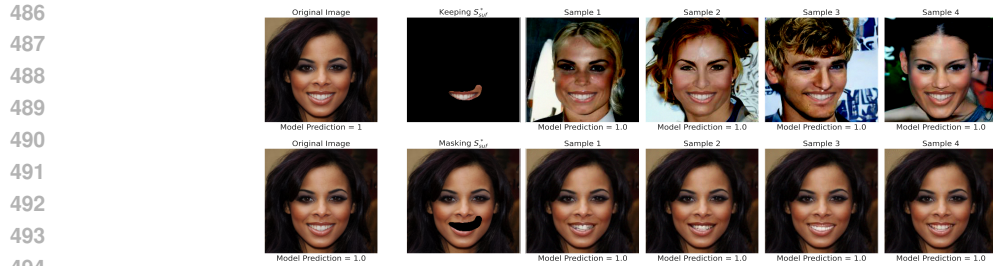
Comparison of Post-hoc Interpretability Methods. For a set of 100 images labeled with a smile and correctly classified by the ResNet classifier, we apply multiple post-hoc interpretability methods and our sufficient and necessary explainers to identify important features associated with smiling. The results in Fig. 3 illustrate that for a wide range of thresholds $t \in [0, 1]$, many methods identify sufficient subsets, as $-\log(\Delta^{\text{suf}})$ for many of them is comparable to that of the sufficient explainer. The necessary explainer, in fact, identifies subsets that are more sufficient than those found by the sufficient explainer. The reason is that the sufficient explainer identifies subsets that are, on average, smaller for all $t \in [0, 1]$, while the necessary explainer finds subsets that are constant in size for all $t \in [0, 1]$ but slightly larger since, to be necessary, they must contain more features that provide additional information about the label. For other methods, as t increases, subset size decreases, and the sufficiency and necessity of the solutions decline. Meanwhile, the necessary explainer naturally identifies necessary subsets, indicated by large $-\log(\Delta^{\text{nec}})$, whereas other methods fail to do so. In conclusion, many methods can identify sufficient sets, but not necessary ones and directly optimizing for these criterion leads to identifying small, constant-sized subsets across thresholds.

Sufficiency vs. Necessity. In Fig. 4, we see how sufficient subsets alone may overlook important features, while solutions to (P_{uni}) offer deeper insights. As stated earlier, the sufficient explainer identifies sets that are sufficient but not necessary. On the other hand, the necessary explainer has high $-\log(\Delta^{\text{suf}})$ and $-\log(\Delta^{\text{nec}})$, indicating that it identifies sufficient *and* necessary set, meaning they also serve as solutions to (P_{uni}) . In Fig. 4, we visualize the reasons for this phenomena. Notice that S_{suf}^* precisely highlights (only) the smile. When S_{suf}^* is fixed, one can generate new images (as done in (Zhang et al., 2023)) for which the model produces the same predictions as it did for the original image (a smile). On the other hand, we also see why S_{suf}^* is *not* necessary: we can fix the complement $(S_{\text{suf}}^*)^c$ and, since there are important features in it, a smile is consistently generated, and the model produces the same prediction on these images as it did on the original. Conversely solutions to (P_{nec}) (also solutions to (P_{uni}) here) generate different explanations that provide a more complete picture of feature importance. Notice that S_{nec}^* is sufficient because $S_{\text{suf}}^* \subseteq S_{\text{nec}}^*$ with the additional features mainly being the dimples and eyes, which aid in determining the presence of a smile. More importantly, Fig. 5 illustrates why S_{nec}^* is necessary: when we fix the complement of S_{nec}^* and generate new samples, half of the faces lack a smile, leading the model f to predict no smile. Additional images and details on sample generation are in Appendices A.2 and A.4.

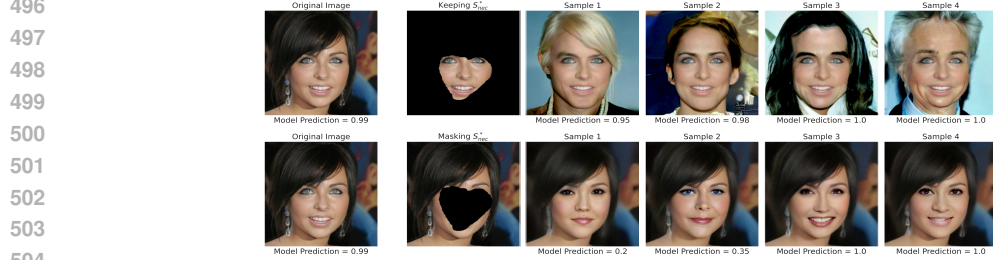
8 LIMITATIONS & BROADER IMPACTS

While this work provides a novel theoretical contribution to the XAI community, there are some limitations that require careful discussion. The choice of reference distribution \mathcal{V}_S determines the characteristics of sufficient and necessary explanations. For instance, only with the true conditional data distribution can one obtain the conditional independence results that our theory provides. Naturally, there are computational trade-offs that must be carefully studied; the ability to learn and sample from accurate conditional distributions to generate explanations with clear statistical meaning comes with a computational and statistical cost, particularly in high-dimensional settings. Thus, a key direction for future work is to explore the impact of different reference distributions and provide a principled framework for selecting a \mathcal{V}_S that balances practical utility and computational feasibility.

Another relevant question is how well our proposed notions align with human intuition. While we aim to understand which features are sufficient and necessary *for a given predicted model*, these explanations may not always correspond to how humans perceive importance (since model might



495 Figure 4: Images and model predictions by fixing and masking the sufficient subset $S_{\text{suf}}^{C^*}$



505 Figure 5: Images and model predictions by fixing and masking the necessary subset $S_{\text{nec}}^{C^*}$

507 use different features to solve a task). This can be an issue in settings where interpretability is
508 essential for trust and accountability, such as in healthcare. On the one hand, our approach can
509 provide useful insights to further evaluate models (e.g. by verifying if the sufficient and necessary
510 features employed by models correlate with the correct ones as informed by human experts). On
511 the other hand, bridging the gap between our mathematical definitions of sufficiency and necessity
512 and other human notions of importance is an area for further investigation. User studies, along
513 with collaboration with domain experts, will be critical in determining how our formal notions of
514 sufficiency and necessity can be adapted or extended to better meet real-world interpretability needs.

515 Finally, the societal impact of this work warrants discussion. While we offer a rigorous framework to
516 understand model predictions, these are oblivious to notions of demographic bias (Hardt et al., 2016;
517 Feldman et al., 2015; Bharti et al., 2024). There is a risk that an “incorrect” choice of generating
518 a sufficient vs. necessary explanation could reinforce biases or obscure the causal reasons behind
519 predictions. Future work will study when and how our framework can incorporate these biases.

520 9 CONCLUSION

521
522 This work formalizes notions of sufficiency and necessity as tools to evaluate feature importance
523 and explain model predictions. We demonstrate that sufficient and necessary explanations, while
524 insightful, often provide incomplete while complementary answers to model behavior. To address
525 this limitation, we propose a unified approach that offers a new and more nuanced understanding
526 of model behavior. Our unified approach expands the scope of explanations and reveals trade-offs
527 between sufficiency and necessity, giving rise to new interpretations of feature importance. Through
528 our theoretical contributions, we present conditions under which sufficiency and necessity align or
529 diverge, and provide two perspectives of our unified approach through the lens of conditional inde-
530 pendence and Shapley values. Our experimental results support our theoretical findings, providing
531 examples of how adjusting sufficiency-necessity trade-off via our unified approach can uncover
532 alternative sets of important features that would be missed by focusing solely on sufficiency or ne-
533 cessity. Furthermore, we evaluate common post-hoc interpretability methods showing that many fail
534 to reliably identify features that are necessary or sufficient. In summary, our work contributes to a
535 more complete understanding of feature importance through sufficiency and necessity. We believe,
536 and hope, our framework holds potential for advancing the rigorous interpretability of ML models.

537 REFERENCES

538 Usha Bhalla, Suraj Srinivas, and Himabindu Lakkaraju. Verifiable feature attributions: A bridge
539 between post hoc explainability and inherent interpretability. *Advances in neural information*

- 540 *processing systems*, 2023.
- 541
- 542 Beepul Bharti, Paul Yi, and Jeremias Sulam. Estimating and controlling for equalized odds via
543 sensitive attribute predictors. *Advances in neural information processing systems*, 36, 2024.
- 544
- 545 Hugh Chen, Ian C Covert, Scott M Lundberg, and Su-In Lee. Algorithms to estimate shapley value
546 feature attributions. *Nature Machine Intelligence*, pp. 1–12, 2023.
- 547
- 548 Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An
549 information-theoretic perspective on model interpretation. In *International conference on ma-*
chine learning, pp. 883–892. PMLR, 2018a.
- 550
- 551 Jianbo Chen, Le Song, Martin J Wainwright, and Michael I Jordan. L-shapley and c-shapley: Effi-
552 cient model interpretation for structured data. *arXiv preprint arXiv:1808.02610*, 2018b.
- 553
- 554 Ian Covert, Scott M Lundberg, and Su-In Lee. Understanding global feature contributions with
555 additive importance measures. *Advances in Neural Information Processing Systems*, 33:17212–
17223, 2020.
- 556
- 557 Ian Covert, Scott Lundberg, and Su-In Lee. Explaining by removing: A unified framework for
558 model explanation. *Journal of Machine Learning Research*, 22(209):1–90, 2021.
- 559
- 560 Piotr Dabkowski and Yarin Gal. Real time image saliency for black box classifiers. *Advances in*
neural information processing systems, 30, 2017.
- 561
- 562 Adnan Darwiche and Auguste Hirth. On the reasons behind decisions. In *ECAI 2020*, pp. 712–720.
563 IOS Press, 2020.
- 564
- 565 Adnan Darwiche and Chunxi Ji. On the computation of necessary and sufficient explanations. In
Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pp. 5582–5591, 2022.
- 566
- 567 Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shan-
568 mugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations
569 with pertinent negatives. *Advances in neural information processing systems*, 31, 2018.
- 570
- 571 Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair
machine learning. *Advances in neural information processing systems*, 34:6478–6490, 2021.
- 572
- 573 Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubra-
574 manian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD*
international conference on knowledge discovery and data mining, pp. 259–268, 2015.
- 575
- 576 Adam E Flanders, Luciano M Prevedello, George Shih, Safwan S Halabi, Jayashree Kalpathy-
577 Cramer, Robyn Ball, John T Mongan, Anouk Stein, Felipe C Kitamura, Matthew P Lungren,
578 et al. Construction of a machine learning dataset through collaboration: the rsna 2019 brain ct
579 hemorrhage challenge. *Radiology: Artificial Intelligence*, 2(3):e190211, 2020.
- 580
- 581 Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal per-
582 turbations and smooth masks. In *Proceedings of the IEEE/CVF international conference on com-*
puter vision, pp. 2950–2958, 2019.
- 583
- 584 Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful pertur-
585 bation. In *Proceedings of the IEEE international conference on computer vision*, pp. 3429–3437,
586 2017.
- 587
- 588 Richard W Hamming. Error detecting and error correcting codes. *The Bell system technical journal*,
29(2):147–160, 1950.
- 589
- 590 Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances*
591 *in neural information processing systems*, 29, 2016.
- 592
- 593 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
770–778, 2016.

- 594 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
595 *neural information processing systems*, 33:6840–6851, 2020.
- 596
- 597 Salim I Amoukou and Nicolas Brunel. Consistent sufficient explanations and minimal local rules for
598 explaining the decision of any classifier or regressor. *Advances in Neural Information Processing*
599 *Systems*, 35:8027–8040, 2022.
- 600 Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva. On relating explanations and adversarial
601 examples. *Advances in neural information processing systems*, 32, 2019.
- 602
- 603 Alexey Ignatiev, Nina Narodytska, Nicholas Asher, and Joao Marques-Silva. From contrastive to
604 abductive explanations and back again. In *International Conference of the Italian Association for*
605 *Artificial Intelligence*, pp. 335–355. Springer, 2020.
- 606 Neil Jethani, Mukund Sudarshan, Yindalon Aphinyanaphongs, and Rajesh Ranganath. Have we
607 learned to explain?: How interpretability methods can learn to encode predictions in their inter-
608 pretations. In *International Conference on Artificial Intelligence and Statistics*, pp. 1459–1467.
609 PMLR, 2021.
- 610
- 611 Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam:
612 Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Pro-*
613 *cessing*, 30:5875–5888, 2021.
- 614 Tero Karras. Progressive growing of gans for improved quality, stability, and variation. *arXiv*
615 *preprint arXiv:1710.10196*, 2017.
- 616
- 617 Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan
618 Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-
619 Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020. URL
620 <https://arxiv.org/abs/2009.07896>.
- 621 Stefan Kolek, Duc Anh Nguyen, Ron Levie, Joan Bruna, and Gitta Kutyniok. A rate-distortion
622 framework for explaining black-box model decisions. In *International Workshop on Extending*
623 *Explainable AI Beyond Deep Models and Classifiers*, pp. 91–115. Springer, 2022.
- 624 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,
625 2015.
- 626
- 627 Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive
628 facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition*
629 *(CVPR)*, 2020.
- 630 Johannes Linder, Alyssa La Fleur, Zibo Chen, Ajasja Ljubetič, David Baker, Sreeram Kannan, and
631 Georg Seelig. Interpreting neural networks for biological sequences by learning stochastic masks.
632 *Nature machine intelligence*, 4(1):41–54, 2022.
- 633
- 634 Ruishan Liu, Shemra Rizzo, Samuel Whipple, Navdeep Pal, Arturo Lopez Pineda, Michael Lu,
635 Brandon Arnieri, Ying Lu, William Capra, Ryan Copping, et al. Evaluating eligibility criteria of
636 oncology trials using real-world data and ai. *Nature*, 592(7855):629–633, 2021.
- 637 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances*
638 *in neural information processing systems*, 30, 2017.
- 639
- 640 Scott M Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M Prutkin, Bala Nair, Ronit
641 Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations to global
642 understanding with explainable ai for trees. *Nature machine intelligence*, 2(1):56–67, 2020.
- 643 Arturo Moncada-Torres, Marissa C van Maaren, Mathijs P Hendriks, Sabine Siesling, and Gijs
644 Geleijnse. Explainable machine learning can outperform cox regression predictions and provide
645 insights in breast cancer survival. *Scientific Reports*, 11(1):1–13, 2021.
- 646
- 647 Edoardo Mosca, Ferenc Szigeti, Stella Tragianni, Daniel Gallagher, and George Louis Groh. Shap-
based explanation methods: A review for nlp interpretability. In *COLING*, 2022.

- 648 Martin Pawelczyk, Klaus Broelemann, and Gjergji Kasneci. On counterfactual explanations under
649 predictive multiplicity. In *Conference on Uncertainty in Artificial Intelligence*, pp. 809–818.
650 PMLR, 2020.
- 651 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ”why should i trust you?” explaining the
652 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference*
653 *on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- 654 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic
655 explanations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 656 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,
657 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-
658 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,
659 2017.
- 660 Lloyd S Shapley. *Notes on the N-person Game*. Rand Corporation, 1951.
- 661 Andy Shih, Arthur Choi, and Adnan Darwiche. A symbolic approach to explaining bayesian net-
662 work classifiers. *arXiv preprint arXiv:1805.03364*, 2018.
- 663 Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through
664 propagating activation differences. In *International conference on machine learning*, pp. 3145–
665 3153. PMLR, 2017.
- 666 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
667 *Advances in neural information processing systems*, 32, 2019.
- 668 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
669 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-*
670 *ional Conference on Learning Representations*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=PXTIG12RRHS)
671 [forum?id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).
- 672 Mukund Sundararajan and Amir Najmi. The many shapley values for model explanation. In *Inter-*
673 *national conference on machine learning*, pp. 9269–9278. PMLR, 2020.
- 674 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
675 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 676 Wesley Tansey, Victor Veitch, Haoran Zhang, Raul Rabadan, and David M Blei. The holdout ran-
677 domization test for feature selection in black box models. *Journal of Computational and Graph-*
678 *ical Statistics*, 31(1):151–162, 2022.
- 679 Jacopo Teneggi, Alexandre Luster, and Jeremias Sulam. Fast hierarchical games for image explana-
680 tions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4494–4503, 2022.
- 681 Jacopo Teneggi, Beepul Bharti, Yaniv Romano, and Jeremias Sulam. Shap-xrt: The shapley value
682 meets conditional independence testing. *Transactions on Machine Learning Research*, 2023.
- 683 The White House. Executive order on the safe, secure, and trustworthy development and use of
684 artificial intelligence, 2023.
- 685 Richard Tomsett, Dave Braines, Dan Harborne, Alun Preece, and Supriyo Chakraborty. Inter-
686 pretable to whom? a role-based model for analyzing interpretable machine learning systems.
687 *arXiv preprint arXiv:1806.07552*, 2018.
- 688 Eric Wang, Pasha Khosravi, and Guy Van den Broeck. Probabilistic sufficient explanations. *arXiv*
689 *preprint arXiv:2105.10118*, 2021.
- 690 Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,
691 Shengchao Liu, Peter Katwyk, Andreea Deac, Anima Anandkumar, Karianne Bergen, Carla P.
692 Gomes, and Shir. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):
693 47–60, August 2023.

702 David Watson, Joshua O’Hara, Niek Tax, Richard Mudd, and Ido Guy. Explaining predictive un-
703 certainty with information theoretic shapley values. *Advances in Neural Information Processing*
704 *Systems*, 36, 2024.

705
706 David S. Watson, Limor Gultchin, Ankur Taly, and Luciano Floridi. Local explanations via necessity
707 and sufficiency: unifying theory and practice. In Cassio de Campos and Marloes H. Maathuis
708 (eds.), *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*,
709 volume 161 of *Proceedings of Machine Learning Research*, pp. 1382–1392. PMLR, 27–30 Jul
710 2021.

711 Jinsung Yoon, James Jordon, and Mihaela van der Schaar. Invase: Instance-wise variable selection
712 using neural networks. In *International Conference on Learning Representations*, 2018.

713
714 Carlos Zednik. Solving the black box problem: a normative framework for explainable artificial
715 intelligence. *Philosophy & Technology*, 34(2):265–288, 2021.

716
717 Guanhua Zhang, Jiabao Ji, Yang Zhang, Mo Yu, Tommi Jaakkola, and Shiyu Chang. Towards
718 coherent image inpainting using denoising diffusion implicit models. In *International Conference*
on Machine Learning, pp. 41164–41193. PMLR, 2023.

719
720 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
721 features for discriminative localization. In *Proceedings of the IEEE conference on computer*
722 *vision and pattern recognition*, pp. 2921–2929, 2016.

723
724 Yazeed Zoabi, Shira Deri-Rozov, and Noam Shomron. Machine learning-based prediction of covid-
725 19 diagnosis based on symptoms. *npj digital medicine*, 4(1):1–5, 2021.

726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 A APPENDIX

757 A.1 PROOFS

759 A.1.1 PROOF OF LEMMA 4.1

760 **Lemma 4.1.** Let $\alpha \in (0, 1)$. For $\tau > 0$, denote S^* to be a solution to (P_{uni}) for which
761 $\Delta_{\mathcal{V}}^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) = \epsilon$. Then, S^* is $\frac{\epsilon}{\alpha}$ -sufficient and $\frac{\epsilon}{1-\alpha}$ -necessary. Formally,

$$762 \quad 0 \leq \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{\alpha} \quad \text{and} \quad 0 \leq \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{1-\alpha}. \quad (9)$$

763 *Proof.* Let $\tau > 0$ and $\alpha \in (0, 1)$ and denote S^* to be a solution to (P_{uni}) such that

$$764 \quad \Delta_{\mathcal{V}}^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) = \epsilon. \quad (10)$$

765 Then, by definition of being a solution to (P_{uni}) ,

$$766 \quad |S^*| \leq \tau. \quad (11)$$

767 Furthermore, recall that

$$768 \quad \Delta_{\mathcal{V}}^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) = \alpha \cdot \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) + (1 - \alpha) \cdot \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \quad (12)$$

769 which implies

$$770 \quad \alpha \cdot \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) = \epsilon - (1 - \alpha) \cdot \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \quad (13)$$

$$771 \quad \leq \epsilon \quad ((1 - \alpha), \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \geq 0) \quad (14)$$

$$772 \quad \implies \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{\alpha}. \quad (15)$$

773 Similarly,

$$774 \quad (1 - \alpha) \cdot \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) = \epsilon - \alpha \cdot \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) \quad (16)$$

$$775 \quad \leq \epsilon \quad (\alpha, \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) \geq 0) \quad (17)$$

$$776 \quad \implies \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \leq \frac{\epsilon}{1-\alpha}. \quad (18)$$

777 \square

778 A.1.2 PROOF OF LEMMA 4.2

779 **Lemma 4.2.** For $0 \leq \epsilon < \frac{\rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x}))}{2}$, denote S_{suf}^* and S_{nec}^* to be ϵ -sufficient and ϵ -necessary sets.
780 Then, if S_{suf}^* is ϵ -super sufficient or S_{nec}^* is ϵ -super necessary,

$$781 \quad S_{\text{suf}}^* \cap S_{\text{nec}}^* \neq \emptyset. \quad (19)$$

782 *Proof.* We will prove the result via contradiction. First recall that,

$$783 \quad f_S(\mathbf{x}) = \mathbb{E}_{\mathbf{x}_{S^c} \sim \mathcal{V}_{S^c}} [f(\mathbf{x}_S, \mathbf{X}_{S^c})] \quad (20)$$

784 and, for any metric $\rho : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$,

$$785 \quad \Delta_{\mathcal{V}}^{\text{suf}}(S, f, \mathbf{x}) \triangleq \rho(f(\mathbf{x}), f_S(\mathbf{x})) \quad (21)$$

$$786 \quad \Delta_{\mathcal{V}}^{\text{nec}}(S, f, \mathbf{x}) \triangleq \rho(f_{S^c}(\mathbf{x}), f_{\emptyset}(\mathbf{x})). \quad (22)$$

787 Since ρ is a metric on \mathbb{R} , it satisfies the triangle inequality. Thus, for $a, b, c \in \mathbb{R}$

$$788 \quad \rho(a, c) \leq \rho(a, b) + \rho(b, c). \quad (23)$$

789 Now, let S_{suf}^* be ϵ -super sufficient and suppose

$$790 \quad S_{\text{suf}}^* \cap S_{\text{nec}}^* = \emptyset. \quad (24)$$

This implies

$$S_{\text{suf}}^* \subseteq (S_{\text{nec}}^*)^c. \quad (25)$$

Subsequently, since S_{suf}^* is ϵ -super sufficient,

$$\Delta_{\mathcal{V}}^{\text{suf}}((S_{\text{nec}}^*)^c, f, \mathbf{x}) \leq \epsilon. \quad (26)$$

As a result, observe

$$\rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x})) \leq \rho(f(\mathbf{x}), f_{(S_{\text{nec}}^*)^c}(\mathbf{x})) + \rho(f_{(S_{\text{nec}}^*)^c}(\mathbf{x}), f_{\emptyset}(\mathbf{x})) \quad \text{triangle inequality} \quad (27)$$

$$= \Delta_{\mathcal{V}}^{\text{suf}}((S_{\text{nec}}^*)^c, f, \mathbf{x}) + \Delta_{\mathcal{V}}^{\text{nec}}((S_{\text{nec}}^*)^c, f, \mathbf{x}) \quad (28)$$

$$\leq \epsilon + \Delta_{\mathcal{V}}^{\text{nec}}((S_{\text{nec}}^*)^c, f, \mathbf{x}) \quad S_{\text{suf}}^* \text{ is } \epsilon\text{-super sufficient} \quad (29)$$

$$\leq 2\epsilon \quad S_{\text{nec}}^* \text{ is } \epsilon\text{-necessary} \quad (30)$$

$$\implies \epsilon \geq \frac{\rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x}))}{2} \quad (31)$$

which is a contradiction because $0 \leq \epsilon < \frac{\rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x}))}{2}$. Thus $S_{\text{suf}}^* \cap S_{\text{nec}}^* \neq \emptyset$. The proof of this result assuming S_{nec}^* is ϵ -super necessary follows the same argument. \square

A.1.3 PROOF OF THEOREM 4.1

Theorem 4.1. Let $\tau_1, \tau_2 > 0$ and $0 \leq \epsilon < \frac{1}{2} \cdot \rho(f(\mathbf{x}), f_{\emptyset}(\mathbf{x}))$. Denote S_{suf}^* and S_{nec}^* to be ϵ -super sufficient and ϵ -super necessary solutions to (P_{suf}) and (P_{nec}) , respectively, such that $|S_{\text{suf}}^*| = \tau_1$ and $|S_{\text{nec}}^*| = \tau_2$. Then, there exists a set S^* such that

$$\Delta_{\mathcal{V}}^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) \leq \epsilon \quad \text{and} \quad \max(\tau_1, \tau_2) \leq |S^*| < \tau_1 + \tau_2. \quad (32)$$

Furthermore, if $S_{\text{suf}}^* \subseteq S_{\text{nec}}^*$ or $S_{\text{nec}}^* \subseteq S_{\text{suf}}^*$, then $S^* = S_{\text{nec}}^*$ or $S^* = S_{\text{suf}}^*$, respectively.

Proof. Consider the set $S^* = S_{\text{suf}}^* \cup S_{\text{nec}}^*$. This set has the following properties:

(P1) S^* is ϵ -sufficient because S_{suf}^* is ϵ -super sufficient

(P2) S^* is ϵ -necessary because S_{suf}^* is ϵ -super necessary

(P3) $|S^*| \geq \max(\tau_1, \tau_2)$ with $|S^*| = \tau_1$ when $S_{\text{nec}}^* \subset S_{\text{suf}}^*$ and with $|S^*| = \tau_2$ when $S_{\text{suf}}^* \subset S_{\text{nec}}^*$

(P4) Via Lemma 4.1, we know $S_{\text{suf}}^* \cap S_{\text{nec}}^* \neq \emptyset$ thus $|S^*| < \tau_1 + \tau_2$

Then by (P1) and (P2)

$$\Delta_{\mathcal{V}}^{\text{uni}}(S^*, f, \mathbf{x}, \alpha) = \alpha \cdot \Delta_{\mathcal{V}}^{\text{suf}}(S^*, f, \mathbf{x}) + (1 - \alpha) \cdot \Delta_{\mathcal{V}}^{\text{nec}}(S^*, f, \mathbf{x}) \quad (33)$$

$$\leq \alpha \cdot \epsilon + (1 - \alpha) \cdot \epsilon = \epsilon \quad (34)$$

and by (P3) and (P4) we have $\max(\tau_1, \tau_2) \leq |S^*| < \tau_1 + \tau_2$, \square

A.1.4 PROOF OF COROLLARY 5.1

Corollary 5.1. Suppose for any $S \subseteq [d]$, $\mathcal{V}_S = p(\mathbf{X}_S | \mathbf{X}_{S^c} = \mathbf{x}_{S^c})$. Let $\alpha \in (0, 1)$, $\epsilon \geq 0$, and denote $\rho : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R}$ to be a metric on \mathbb{R} . Furthermore, for $f(\mathbf{X}) = \mathbb{E}[Y | \mathbf{X}]$ and $\tau > 0$, let S^* be a solution to (P_{uni}) such that $\Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) = \epsilon$. Then, S^* satisfies the following conditional independence relations,

$$\rho(\mathbb{E}[Y | \mathbf{x}], \mathbb{E}[Y | \mathbf{X}_{S^*} = \mathbf{x}_{S^*}]) \leq \frac{\epsilon}{\alpha} \quad \text{and} \quad \rho(\mathbb{E}[Y | \mathbf{X}_{S^c} = \mathbf{x}_{S^c}], \mathbb{E}[Y]) \leq \frac{\epsilon}{1 - \alpha}. \quad (35)$$

Proof. All we need to show is that when $\mathcal{V}_S = p(\mathbf{X}_S | \mathbf{X}_{S^c} = \mathbf{x}_{S^c})$ and $f(\mathbf{X}) = \mathbb{E}[Y | \mathbf{X}]$, we have

$$f_S(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X}_S = \mathbf{x}_S]. \quad (36)$$

864 Once this is proven, we can simply apply Lemma 4.1.

865 To this end, we have by assumption that $f(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$ and, for any $S \subseteq [d]$, $\nu_S = p(\mathbf{X}_S \mid$
866 $\mathbf{X}_{S^c} = \mathbf{x}_{S^c})$. Then by definition

$$867 f_S(\mathbf{x}) = \mathbb{E}_{\nu_{S^c}}[f(\mathbf{x}_S, \mathbf{X}_{S^c})] = \int_{\mathcal{X}} f(\mathbf{x}_S, \mathbf{X}_{S^c}) \cdot p(\mathbf{X}_{S^c} \mid \mathbf{X}_S = \mathbf{x}_S) d\mathbf{X}_{S^c} \quad (37)$$

$$870 = \int_{\mathcal{X}} \mathbb{E}[Y \mid \mathbf{X}_S = \mathbf{x}_S, \mathbf{X}_{S^c}] \cdot p(\mathbf{X}_{S^c} \mid \mathbf{X}_S = \mathbf{x}_S) d\mathbf{X}_{S^c} \quad (38)$$

$$871 = \int_{\mathcal{X}} \left(\int_{\mathcal{Y}} y \cdot p(y \mid \mathbf{X}_S = \mathbf{x}_S, \mathbf{X}_{S^c}) dy \right) \cdot p(\mathbf{X}_{S^c} \mid \mathbf{X}_S = \mathbf{x}_S) d\mathbf{X}_{S^c} \quad (39)$$

$$872 = \int_{\mathcal{Y}} y \left(\int_{\mathcal{X}} p(y, \mathbf{X}_{S^c} \mid \mathbf{X}_S = \mathbf{x}_S) d\mathbf{X}_{S^c} \right) dy \quad (40)$$

$$873 = \int_{\mathcal{Y}} y \cdot p(y \mid \mathbf{X}_S = \mathbf{x}_S) dy \quad (41)$$

$$874 = \mathbb{E}[Y \mid \mathbf{X}_S = \mathbf{x}_S]. \quad (42)$$

875 By applying Lemma 4.1, we have the desired result. \square

882 A.1.5 PROOF OF THEOREM 5.1

883 **Theorem 5.1.** Consider an input \mathbf{x} for which $f(\mathbf{x}) \neq f_\emptyset(\mathbf{x})$. Denote by $\Lambda_d = \{S, S^c\}$ the partition
884 of $[d] = \{1, 2, \dots, d\}$, and define the characteristic function to be $v(S) = -\rho(f(\mathbf{x}), f_S(\mathbf{x}))$. Then,

$$885 \phi_S^{\text{shap}}(\Lambda_d, v) \geq \rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha). \quad (43)$$

886 *Proof.* Before we prove the result, recall the following properties of a metric ρ in the reals:

$$887 \text{(P1)} \quad \forall a, b \in \mathbb{R}, \quad \rho(a, b) = 0 \iff a = b$$

$$888 \text{(P2)} \quad \text{for } a, b, c \in \mathbb{R}, \quad \rho(a, c) \leq \rho(a, b) + \rho(b, c).$$

889 Now, for the partition $\Lambda_d = \{S, S^c\}$ of $[d] = \{1, 2, \dots, d\}$ and characteristic function $v(S) =$
890 $-\rho(f(\mathbf{x}), f_S(\mathbf{x}))$, $\phi_S^{\text{shap}}(\Lambda_d, v)$ is defined as

$$891 \phi_S^{\text{shap}}(\Lambda_d, v) = \frac{1}{2} \cdot [v(S \cup S^c) - v(S^c)] + \frac{1}{2} \cdot [v(S) - v(\emptyset)] \quad (44)$$

$$892 = \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_{S^c}(\mathbf{x})) - \rho(f(\mathbf{x}), f(\mathbf{x}))] + \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f(\mathbf{x}), f_S(\mathbf{x}))] \quad (45)$$

$$893 = \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_{S^c}(\mathbf{x}))] + \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f(\mathbf{x}), f_S(\mathbf{x}))] \quad \text{by (P1)} \quad (46)$$

894 By (P2)

$$895 \rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) \leq \rho(f(\mathbf{x}), f_{S^c}(\mathbf{x})) + \rho(f_{S^c}(\mathbf{x}), f_\emptyset(\mathbf{x})) \quad (47)$$

$$896 \implies \rho(f(\mathbf{x}), f_{S^c}(\mathbf{x})) \geq \rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f_{S^c}(\mathbf{x}), f_\emptyset(\mathbf{x})). \quad (48)$$

897 Thus

$$898 \phi_S^{\text{shap}}(\Lambda_d, v) = \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_{S^c}(\mathbf{x}))] + \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f(\mathbf{x}), f_S(\mathbf{x}))] \quad (49)$$

$$899 \geq \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f_{S^c}(\mathbf{x}), f_\emptyset(\mathbf{x}))] + \frac{1}{2} \cdot [\rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \rho(f(\mathbf{x}), f_S(\mathbf{x}))] \quad (50)$$

$$900 = \rho(f(\mathbf{x}), f_\emptyset(\mathbf{x})) - \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha). \quad (51)$$

901 \square

918 A.2 ADDITIONAL EXPERIMENTAL DETAILS

919
920 In this section, we include further experimental details. All experiments were performed on a private
921 cluster with 8 NVIDIA RTX A5000 with 24 GB of memory. All scripts were run on PyTorch
922 2.0.1, Python 3.11.5, and CUDA 12.2.

923 A.2.1 RSNA CT HEMORRHAGE

924
925 **Dataset Details.** The RSNA 2019 Brain CT Hemorrhage Challenge dataset (Flanders et al., 2020),
926 contains 752803 images labeled by a panel of board-certified radiologists with the types of hemor-
927 rhage present (epidural, intraparenchymal, intraventricular, subarachnoid, subdural).
928

929 **Implementation.** Recall for this experiment, to identify sufficient and necessary masks S for a
930 sample \mathbf{x} , we considered the relaxed optimization problem (Fong et al., 2019; Kolek et al., 2022)

$$931 \arg \min_{S \subseteq [0,1]^d} \Delta_{\mathcal{V}}^{\text{uni}}(S, f, \mathbf{x}, \alpha) + \lambda_1 \cdot \|S\|_1 + \lambda_{\text{TV}} \cdot \|S\|_{\text{TV}}. \quad (52)$$

932
933 where $\|S\|_1$ and $\|S\|_{\text{TV}}$ are the L^1 and Total Variation norm of S , which promote sparsity and
934 smoothness respectively and λ_{sp} and λ_{sm} are the associated. To solve this problem, a mask
935 $S \in [0, 1]^{512 \times 512}$ is initialized with entries $S_i \sim \mathcal{N}(0.5, \frac{1}{36})$. For 1000 iterations, the mask S
936 is iteratively updated to minimize
937

$$938 \alpha \cdot |f(\mathbf{x}) - f_S(\mathbf{x})| + (1 - \alpha) \cdot |f(\mathbf{x}) - f_S(\mathbf{x})| + \lambda_1 \cdot \|S\|_1 + \lambda_{\text{TV}} \cdot \|S\|_{\text{TV}} \quad (53)$$

939 where for any S ,

$$940 f_S(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K f((\tilde{\mathbf{X}}_S)_i) \quad \text{with} \quad (\tilde{\mathbf{X}}_S)_i = \mathbf{x} \circ \tilde{\mathbb{1}}_S + (1 - \tilde{\mathbb{1}}_S) \circ b_i. \quad (54)$$

941
942 Here the entries $(\tilde{\mathbb{1}}_S)_i \sim \text{Bernoulli}(S_i)$ and b_i is the i th entry of a vector $\mathbf{b} = (b_1, \dots, b_d) \sim \mathcal{V}$. In
943 our implementation the reference distribution \mathcal{V} is the unconditional mean image over the of training
944 images and so b_i is the simply the average value of the i th pixel over the training set. To allow for
945 differentiation during optimization, we generate discrete samples $\tilde{\mathbb{1}}_S$ using the Gumbel-Softmax
946 distribution. This methodology simply implies the entries $(\tilde{\mathbf{X}}_S)_i$ is a Bernoulli distribution with
947 outcomes $\{b_i, x_i\}$, i.e. $(\tilde{\mathbf{X}}_S)_i$ is distributed as

$$948 \Pr[(\tilde{\mathbf{X}}_S)_i = x_i] = S_i \quad (55)$$

$$949 \Pr[(\tilde{\mathbf{X}}_S)_i = b_i] = 1 - S_i \quad (56)$$

950
951 For each $\alpha \in \{0, 0.5, 1\}$, during optimization we set $K = 10$, $\lambda_1 = 3$ and $\lambda_{\text{TV}} = 20$ and use the
952 Adam optimizer with default β -parameters of $\beta_1 = 0.9$, $\beta_2 = 0.99$ and a fixed learning rate of 0.01.
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

A.2.2 CELEBA-HQ

Dataset Details. We use a modified version of the CelebA-HQ dataset (Lee et al., 2020; Karras, 2017) which contains 30,000 celebrity faces resized to 256×256 pixels with several landmark locations and binary attributes (e.g., eyeglasses, bangs, smiling).

Implementation. Recall for this experiment, to generate sufficient or necessary masks S for samples \mathbf{x} , we learn a model $g_\theta : \mathcal{X} \mapsto [0, 1]^d$ via solving the following optimization problem:

$$\arg \min_{\theta \in \Theta} \mathbb{E}_{\mathbf{X} \sim \mathcal{D}_X} \left[\Delta_{\mathcal{V}}^{\text{uni}}(g_\theta(\mathbf{X}), f, \mathbf{X}, \alpha) + \lambda_1 \cdot \|g_\theta(\mathbf{X})\|_1 + \lambda_{\text{TV}} \cdot \|g_\theta(\mathbf{X})\|_{\text{TV}} \right] \quad (57)$$

To learn sufficient and necessary explainer models, we solve Eq. (8) via empirical risk minimization for $\alpha \in \{0, 1\}$ respectively. Given N samples $\{\mathbf{X}_i\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_X$, we solve

$$\frac{1}{N} \sum_{i=1}^N \left[\Delta_{\mathcal{V}}^{\text{uni}}(g_\theta(\mathbf{X}_i), f, \mathbf{X}_i, \alpha) + \lambda_1 \cdot \|g_\theta(\mathbf{X}_i)\|_1 + \lambda_{\text{TV}} \cdot \|g_\theta(\mathbf{X}_i)\|_{\text{TV}} \right]. \quad (58)$$

Here

$$\Delta_{\mathcal{V}}^{\text{uni}}(g_\theta(\mathbf{x}_i), f, \mathbf{x}_i, \alpha) = \alpha \cdot |f(\mathbf{x}_i) - f_S(\mathbf{x}_i)| + (1 - \alpha) \cdot |f(\mathbf{x}_i) - f_S(\mathbf{x}_i)| \quad (59)$$

where $f_S(\mathbf{x}_i)$ is evaluated in the same manner as in the RSNA experiment. For $\alpha = 0$, $\lambda_1 = 0.1$ and $\lambda_{\text{TV}} = 100$. For $\alpha = 1$, $\lambda_1 = 1$ and $\lambda_{\text{TV}} = 10$. For both α , during optimization we use a batch size of 32, set $K = 10$ and use the Adam optimizer with default β -parameters of $\beta_1 = 0.9$, $\beta_2 = 0.99$ and a fixed learning rate of 1×10^{-4}

Sampling. To generate the samples in Figs. 4 and 5, samples we use the CoPaint method (Zhang et al., 2023). We utilize their code base and pretrained diffusion models with the exact the same parameters as reported in the paper to perform conditional generation. Everything used is available at <https://github.com/UCSB-NLP-Chang/CoPaint>.

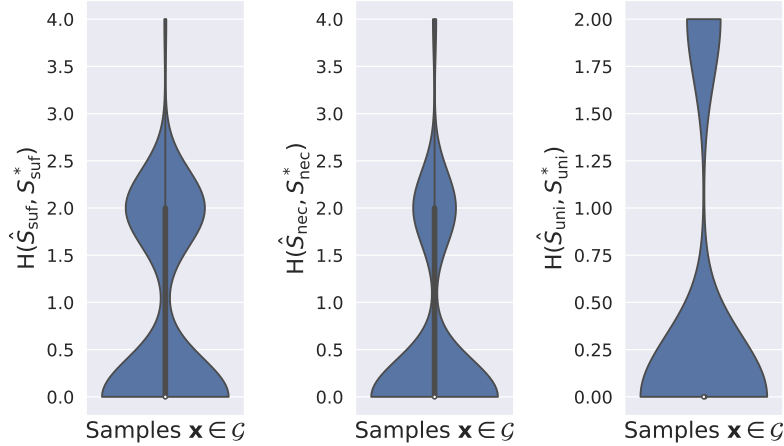


Figure 7: Hamming distances between computed and optimal solutions for P_{suf} , P_{nec} , and P_{uni}

A.3 ADDITIONAL EXPERIMENTS

A.3.1 SYNTHETIC EXAMPLE

We model features $\mathbf{X} \in \mathbb{R}^7$, where $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ for $i \in \{1, 4, 5, 6, 7\}$. The remaining features and response Y follow:

$$X_2 = 2 \cdot X_1 + \epsilon, \quad Y = 4 \cdot X_2 \cdot \mathbf{1}_{\{X_2 > 10\}} + \epsilon, \quad X_3 = 4 \cdot Y + 15 \cdot X_4 \cdot \mathbf{1}_{\{X_4 > 0.5\}} + \epsilon \quad (60)$$

where $\epsilon \sim \mathcal{N}(0, 1)$. For $\mathbf{X} \in \mathcal{G} := \{\mathbf{X} \mid X_2 > 10, X_4 > \frac{1}{2}\}$, the data-generating process is represented by the directed acyclic graph (DAG) shown in Fig. 6 (note X_5, X_6 and X_7 are not depicted since they share no dependencies with any of the random variables). We can see that $Y \perp\!\!\!\perp \mathbf{X}_{\{1,5,6,7\}} \mid \mathbf{X}_{2,3,4}$ and $Y \perp\!\!\!\perp \mathbf{X}_{\{4,5,6,7\}}$.

Thus, for $f(\mathbf{X}) = \mathbb{E}[Y \mid \mathbf{X}]$ and $\mathcal{V}_S = p(\mathbf{X}_{S^c} \mid \mathbf{x}_S)$, the solutions to P_{suf} , P_{nec} , and P_{uni} with $\tau = 4$ are:

$$S_{\text{suf}}^* = \{2, 3, 4\}, \quad S_{\text{nec}}^* = \{1, 2, 3\}, \quad S_{\text{uni}}^* = \{1, 2, 3, 4\}.$$

In this experiment, we train a general predictor (a three-layer fully-connected neural network) to approximate $\mathbb{E}[Y \mid \mathbf{X}]$ and

1. Validate the sets listed above are the optimal solutions.
2. Demonstrate that common post-hoc interpretability methods struggle to recover these solutions.

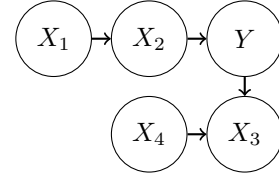


Figure 6: DAG modeling the data-generating process for $\mathbf{X} \in \mathcal{G}$

Validation of Solutions. For type $\in \{\text{suf}, \text{nec}, \text{uni}\}$, $\tau = 4$, and 100 samples $\mathbf{x} \in \mathcal{G}$ we compute solutions to P_{type} , denoted as \hat{S}_{type} , via exhaustive search. Fig. 7 shows that for all three problems, the Hamming distance between \hat{S}_{type} and S_{type}^* is equal to 0 for a majority of the samples in \mathcal{G} . These results indicate that the solutions computed via an exhaustive search do typically retrieve the correct solutions (the minor discrepancies are due to $f(\mathbf{X})$ being an approximation of $\mathbb{E}[Y \mid \mathbf{X}]$). More importantly, this setting is a clear example of how the unified approach provides a different perspective of importance. One would not be able to identify the set $S = \{1, 2, 3, 4\}$ as the most important one without directly solving the unified problem.

Comparison with Post-hoc Methods For our model f and samples $\mathbf{x} \in \mathcal{G}$, we use Integrated Gradients, Gradient Shapley, DeepLift, and Lime to generate attribution scores. To identify whether these methods highlight sufficient and/or necessary features, and as done with our other experiments, we perform the following steps on the attribution scores for a sample \mathbf{x} (so that the outputs of all methods are comparable)

1. We normalize the scores to the interval $[0, 1]$ via min/max normalization.

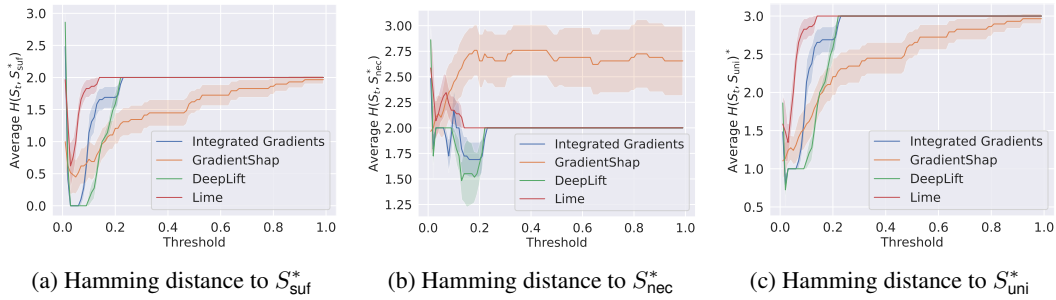


Figure 8: Comparison of various post-hoc methods

2. We generate binary masks S_t by thresholding the normalized scores with thresholds $t \in (0, 1)$
3. For $\text{type} \in \{\text{suf}, \text{nec}, \text{uni}\}$, we compute $H(S_t, S_{\text{type}}^*)$, the Hamming distance between S_t and the true solutions to P_{suf} , P_{nec} , and P_{uni}

The results in Fig. 8 illustrate that, in general, current post-hoc methods fail to recover the optimal sufficient, necessary, or unified solutions. For thresholds $t \in [0, 0.1]$, we see that Integrated Gradients and DeepLift recover solutions S_t that match the optimal sufficient solution S_{suf}^* . This indicates these methods are capable of highlighting the sufficient features. Besides this observation, we see that for thresholds $t > 0.2$ and all three problems, nearly all methods recover solutions S_t that have a Hamming distance ≥ 2 to the optimal solution indicating that the solutions S_t and optimal solutions S^* differ by at least two elements. As a result, the conclusion is that most common methods do *not* detect sufficient solutions and *no* methods detect necessary or unified solutions.

A.4 ADDITIONAL FIGURES

A.4.1 RSNA CT HEMORRHAGE

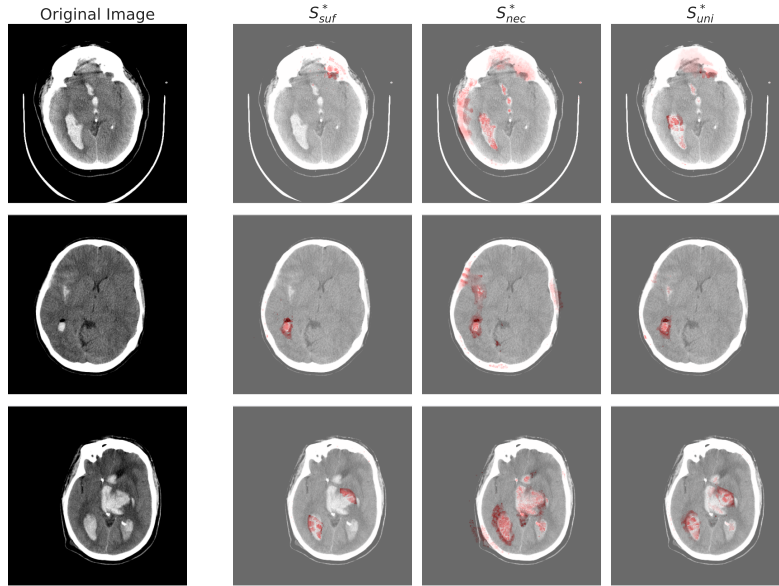


Figure 9: S_{suf}^* , S_{nec}^* and S_{uni}^* for various CT scans.

A.4.2 CELEBA-HQ

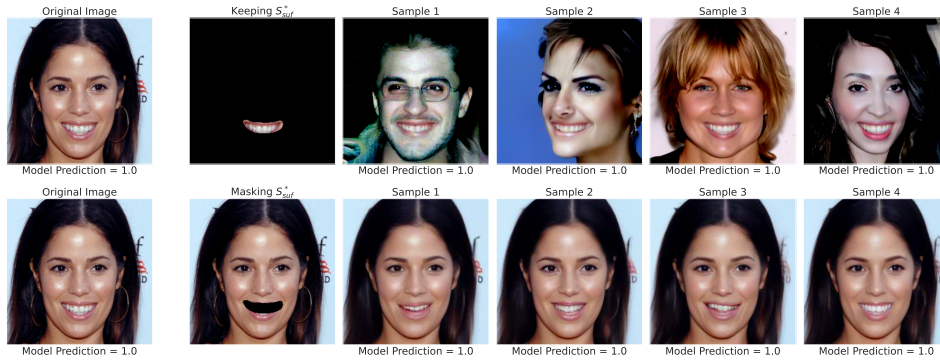


Figure 10: Images and model predictions by fixing and masking the sufficient subset S_{suf}^*

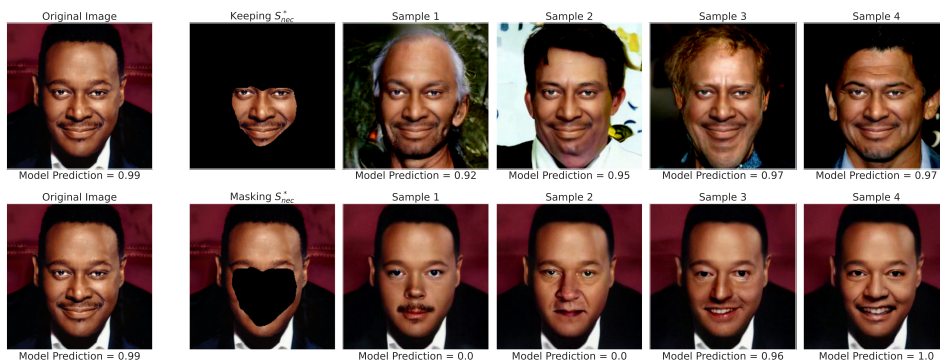


Figure 11: Images and model predictions by fixing and masking the necessary subset S_{nec}^*