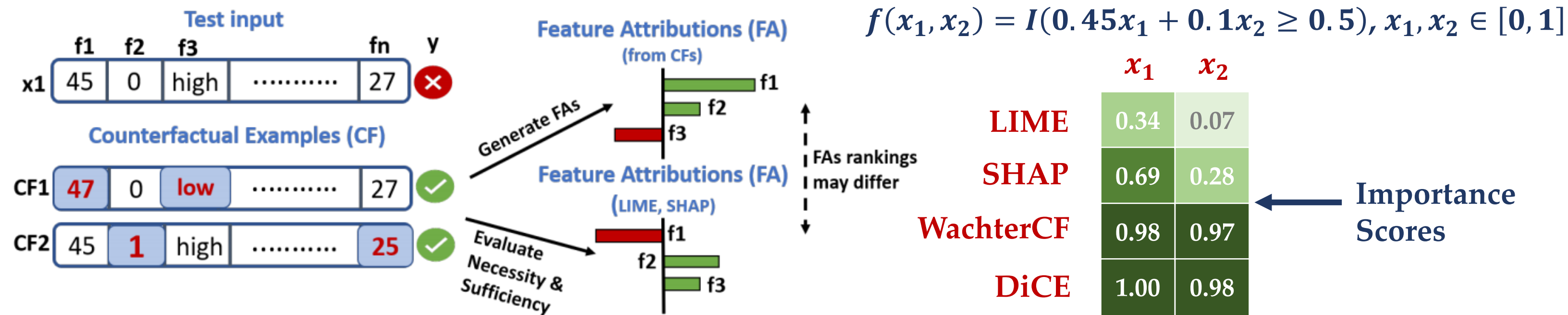


Towards Unifying Feature Attribution and Counterfactual Explanations: Different Means to the Same End

Local Explanation Methods DISAGREE with Each Other

Feature Attributions and Counterfactuals often disagree even for simple linear models



- Propose an **unifying framework** based on **Actual Causality** to interpret these two approaches
- Evaluate** attribution-based methods on the **necessity** and **sufficiency** of their top-ranked features

Actual Causality and Sufficiency → Ideal Model Explanations

- Existence:** There exists a context $u \in U$ such that $x_j = a$ and $f(x_{-j} = b, x_j = a) = y^*$.
- Necessity:** For each context $u \in U$ where $x_j = a$ and $f(x_{-j} = b, x_j = a) = y^*$, some feature subset $x_{sub} \subseteq x_j$ is an actual cause under (M, u)
- Minimality:** x_j is minimal, namely, there is no strict subset $x_s \subset x_j$ such that $x_s = a_s$ satisfies conditions 1-2 above, where $a_s \subset a$.
- Sufficiency:** For all contexts $u' \in U$, $x_j \leftarrow a \Rightarrow y = y^*$.

Stronger **Necessity** condition
(But-for):

Changing the value of x_j **alone** changes the prediction of the model (that is when all other features are kept the same)

Ideal Model Explanations → Partial Model Explanations

- However, for most realistic ML models, an **ideal explanation is impractical**.
 - It is rare to find such clean explanations of a ML model's output
 - Example:** there is **no sufficient feature** for $f(x_1, x_2, x_3) = I(0.4x_1 + 0.1x_2 + 0.1x_3 \geq 0.5)$
- (α, β) **goodness of an explanation** to capture the *extent* to which a feature is necessary or sufficient to "cause" the model's original output

$$\alpha = \Pr(x_j \text{ is a cause of } y^* | x_j = a, y = y^*)$$

$$\beta = \Pr(y = y^* | x_j \leftarrow a)$$

Interpretation Using A Unifying Framework

Counterfactual explanation (α_{CF})

- Optimizes **Necessity**
 - Perturbed feature subset x_j is a **but-for cause** of the original output
 - α_{CF} summarizes the outcomes of all such perturbations and ranks any feature subset for their necessity
- $$\alpha_{CF} = \Pr((x_j \leftarrow a' \Rightarrow y \neq y^*) | x_j = a, y = y^*)$$

Attribution-based explanations (β)

- Optimizes **Sufficiency**
- Importance of x_j can be interpreted as its sufficiency
- The fraction of all contexts where $x_j \leftarrow a$ leads to $y = y^*$ is given by

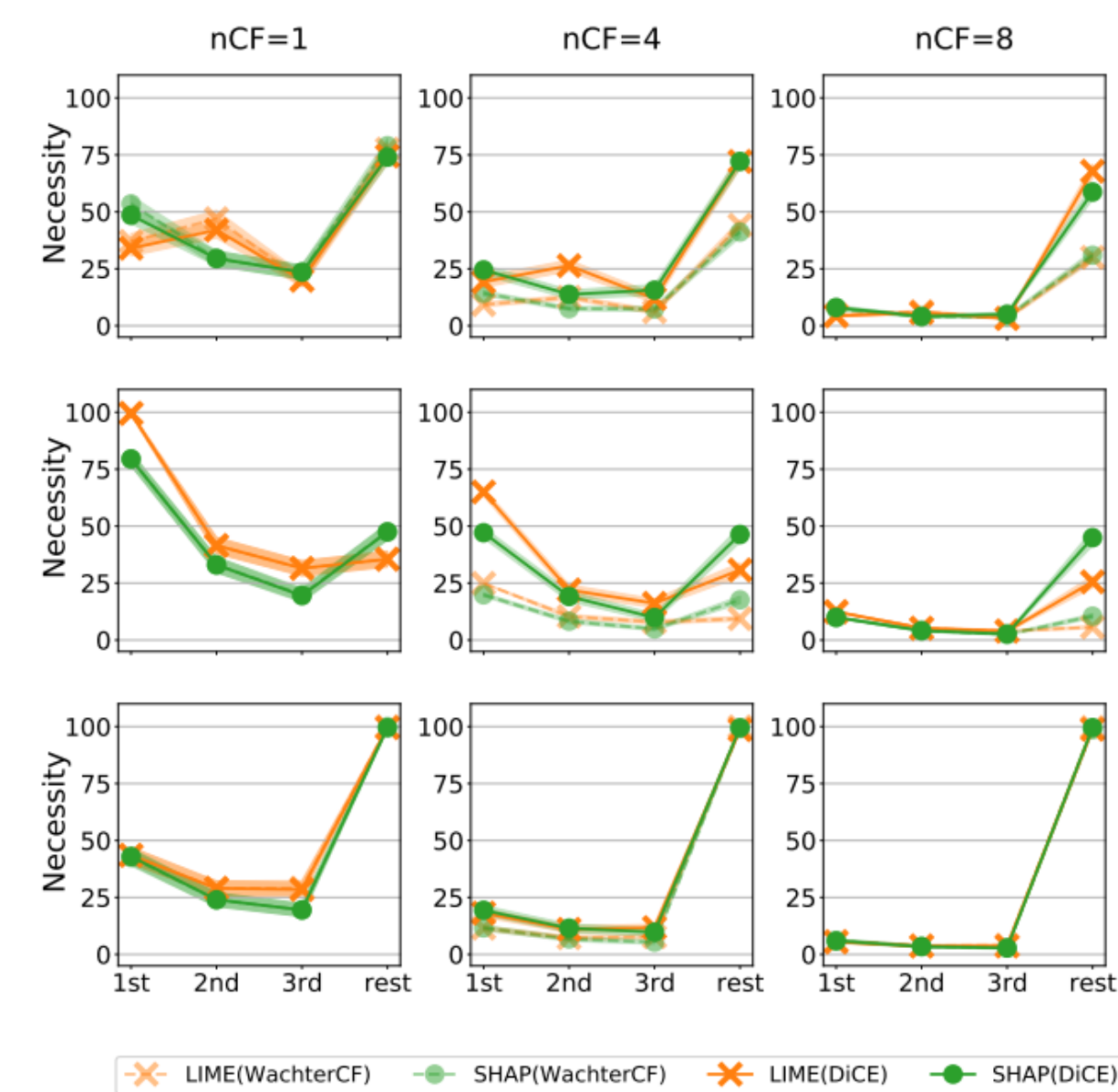
$$\beta = \Pr(y = y^* | x_j \leftarrow a)$$

Top Features of LIME/SHAP are Neither Necessary Nor Sufficient

We use counterfactual explanations to evaluate feature attribution methods based on Necessity and Sufficiency

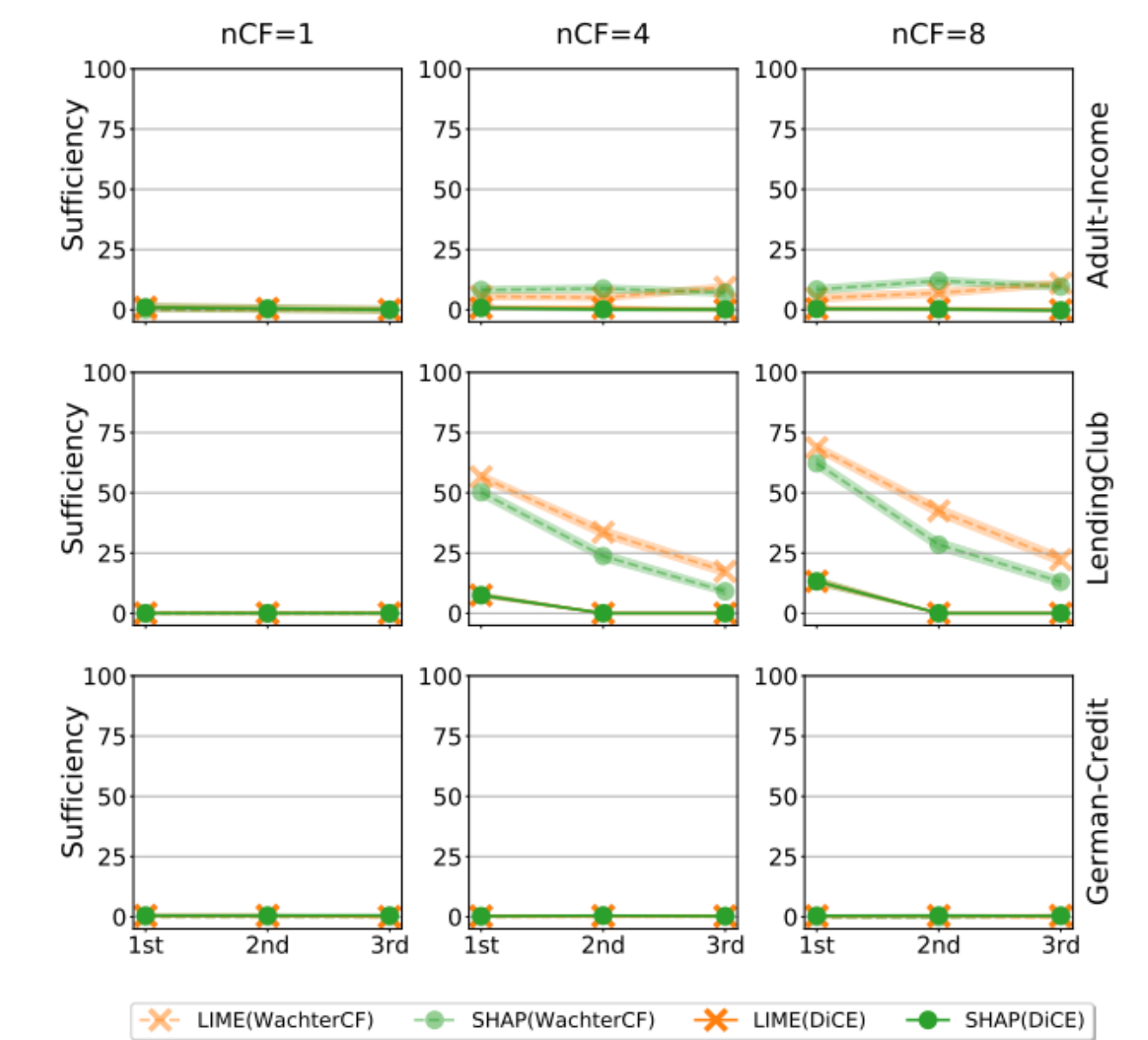
Generate CFs by **changing only** x_j

$$\text{Necessity} = \frac{\sum_{i, x_j \neq a} \mathbb{1}(CF_i)}{nCF * N}$$



Generate CFs by **fixing only** x_j

$$\text{Sufficiency} = \frac{\sum_i \mathbb{1}(CF_i)}{nCF * N} - \frac{\sum_{i, x_j \leftarrow a} \mathbb{1}(CF_i)}{nCF * N}$$



- Highly ranked features may often **neither be necessary nor sufficient** explanations of a model's predictions – Other features are (sometimes more) meaningful and can potentially provide actionable changes
- Necessity and Sufficiency become weaker for top-ranked features as the **number of features** in a dataset **increases**
- Important to consider **multiple explanation methods** to understand the predictions of a ML model