

# Evaluating explanations

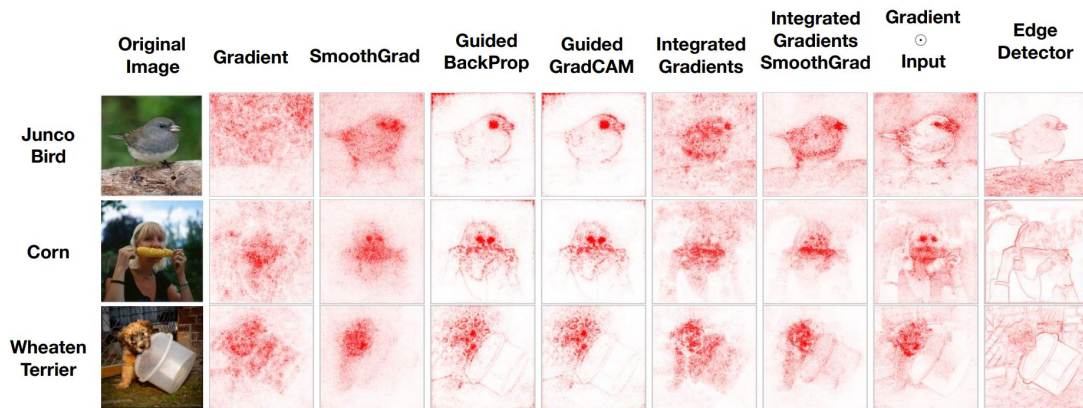
---

Hubert Baniecki

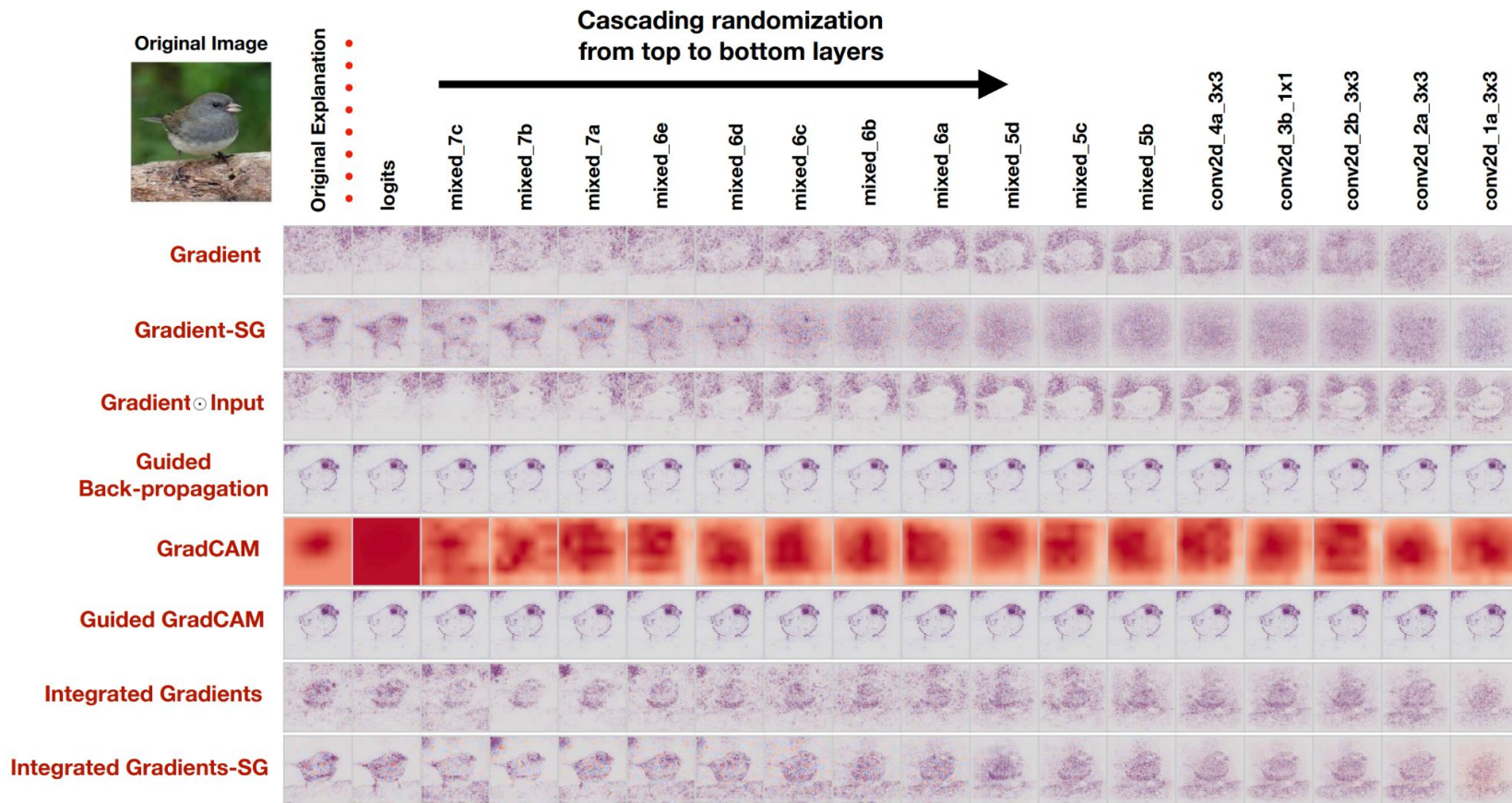
January 10th, 2022

# Motivation testing explanations

“Through extensive experiments we show that some existing saliency methods are independent both of the **model** and of the **data** generating process.”



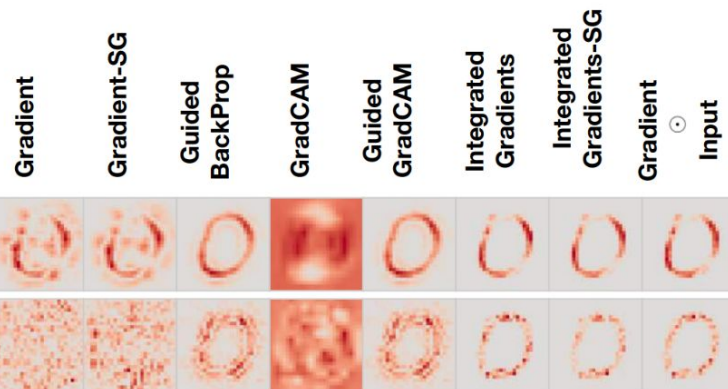
Adebayo et al. Sanity Checks for Saliency Maps (NeurIPS 2018)



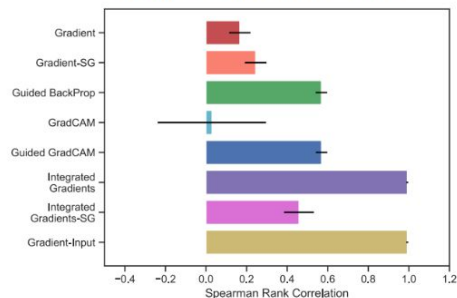
Adebayo et al. Sanity Checks for Saliency Maps (NeurIPS 2018)

## CNN - MNIST

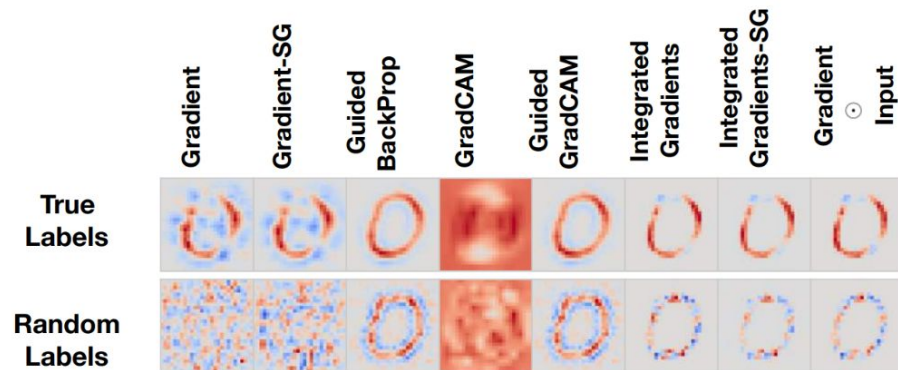
### Absolute-Value Visualization



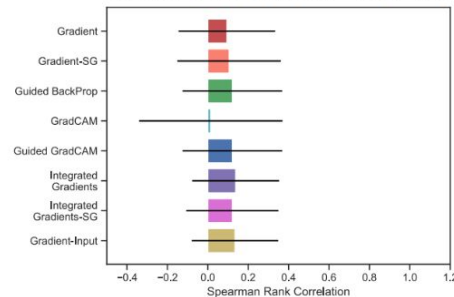
Rank Correlation - Abs



### Diverging Visualization



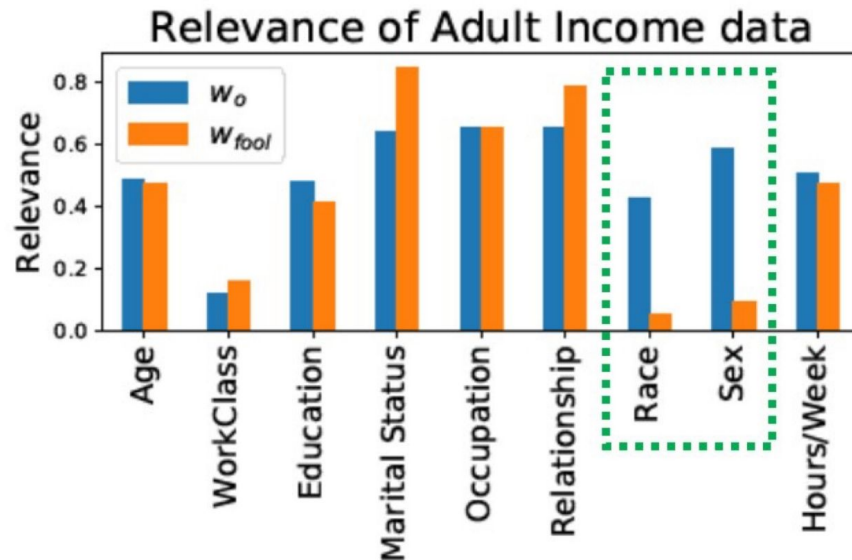
Rank Correlation - No Abs



Adebayo et al. Sanity Checks for Saliency Maps (NeurIPS 2018)



# Motivation manipulating explanations



Dombrowski et al. Explanations can be manipulated and geometry is to blame (NeurIPS 2019)

Heo et al. Fooling Neural Network Interpretations via Adversarial Model Manipulation (NeurIPS 2019)

# Agenda

discuss three ways of evaluating explanations

1. **User** study: Poursabzi-Sangdeh et al. Manipulating and Measuring Model Interpretability (CHI 2021)
2. **Theoretical** benchmark: Liu et al. Synthetic Benchmarks for Scientific Research in Explainable Machine Learning (NeurIPS Dataset Track 2021)
3. **Practical** experiment: Zhou et al. Do Feature Attribution Methods Correctly Attribute Features? (AAAI 2022)

# User study

Poursabzi-Sangdeh et al.

\*Microsoft Research

Manipulating and Measuring  
Model Interpretability  
(CHI 2021)

---

**Idea.** Compare usefulness of four models in the experiments: white-box and black-box with 2 or 8 features used.

- (1) **How well can people simulate a model's predictions?**
- (2) **To what extent do people follow a model's predictions when it is beneficial for them to do so?**
- (3) **How well can people detect when a model has made a mistake and correct for it?**

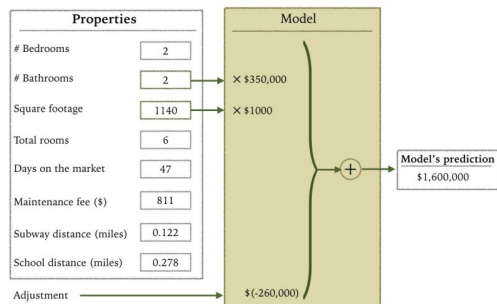
**Methods.** Four sequential user-studies on Amazon Technical Turk.

EXPERIMENT 1: PREDICTING APARTMENT SELLING PRICES (N=1250)

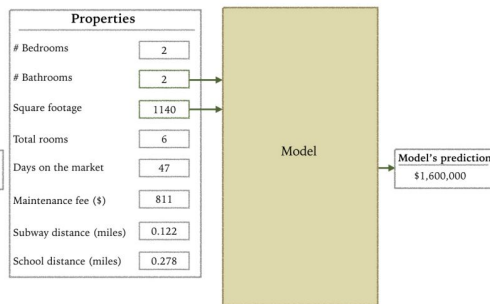
EXPERIMENT 4: OUTLIER FOCUS AND DETECTION OF MISTAKES (N=800)



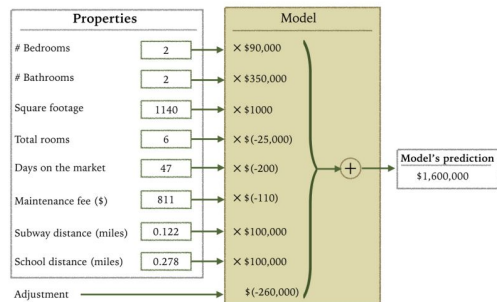
# EXPERIMENT 1: PREDICTING APARTMENT SELLING PRICES (N=1250)



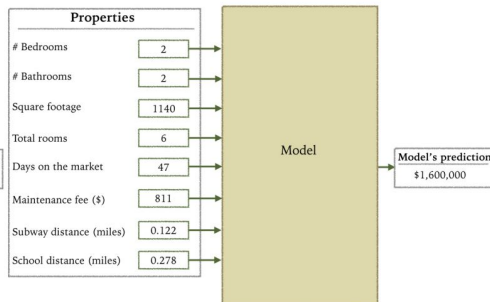
(a) Clear, two-feature condition (CLEAR-2).



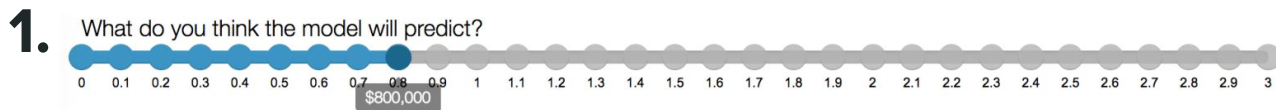
(b) Black-box, two-feature condition (BB-2).



(c) Clear, eight-feature condition (CLEAR-8).



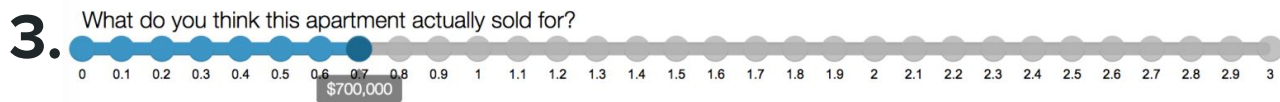
(d) Black-box, eight-feature condition (BB-8).



How confident are you the model will predict this?



How confident are you that the model got it right?

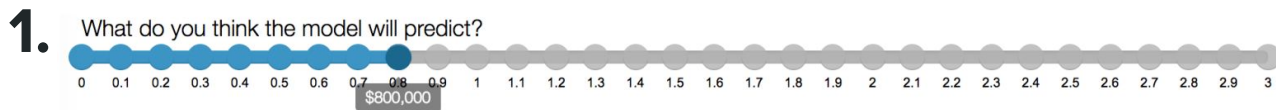


How confident are you that you got it right?



# Hypotheses and measures

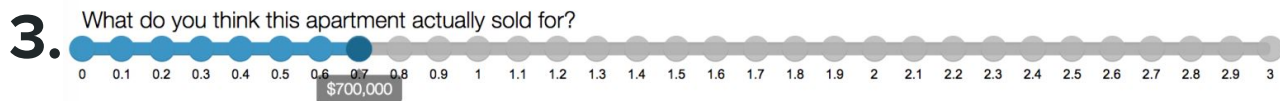
1. **Simulation:** users will better simulate predictions of a simple model  
error:  $\text{model prediction} - \text{user prediction of the prediction}$
2. **Deviation:** users will follow the predictions of a simple model  
error:  $\text{model prediction} - \text{user prediction of the price}$
3. **Detection of mistakes:** users will differently find the models' mistakes  
same as above but only for 2 apartments



How confident are you the model will predict this?



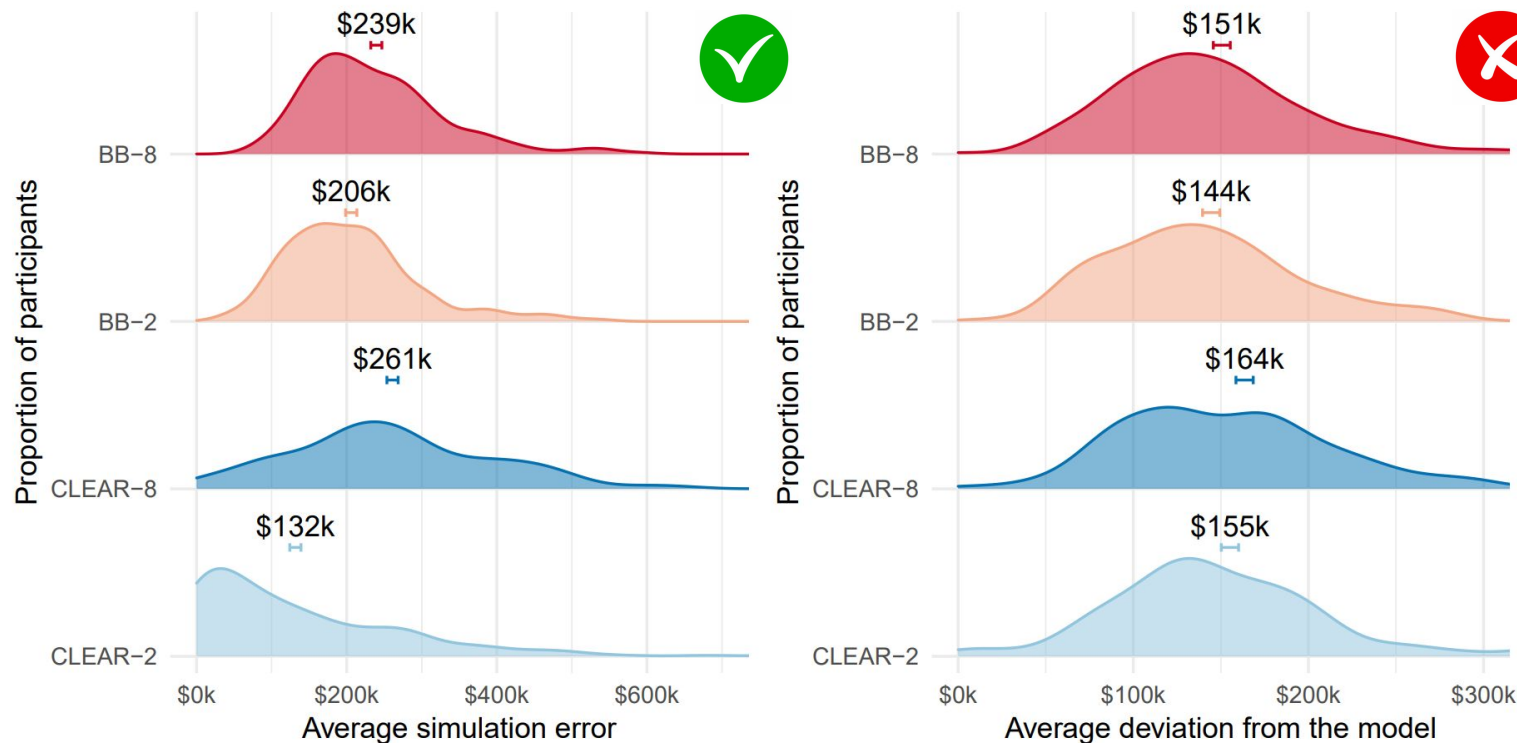
How confident are you that the model got it right?

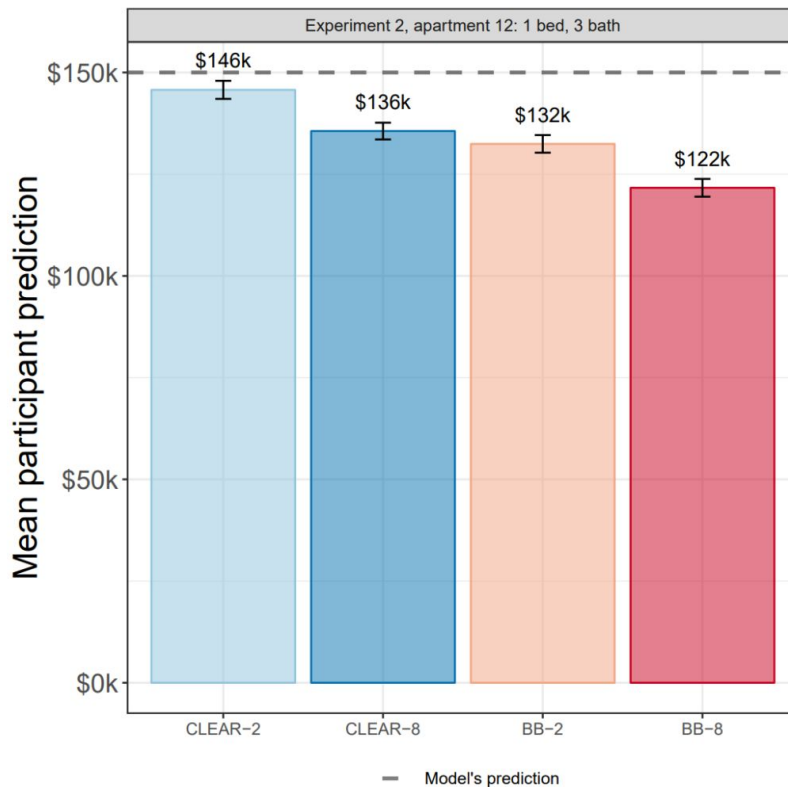
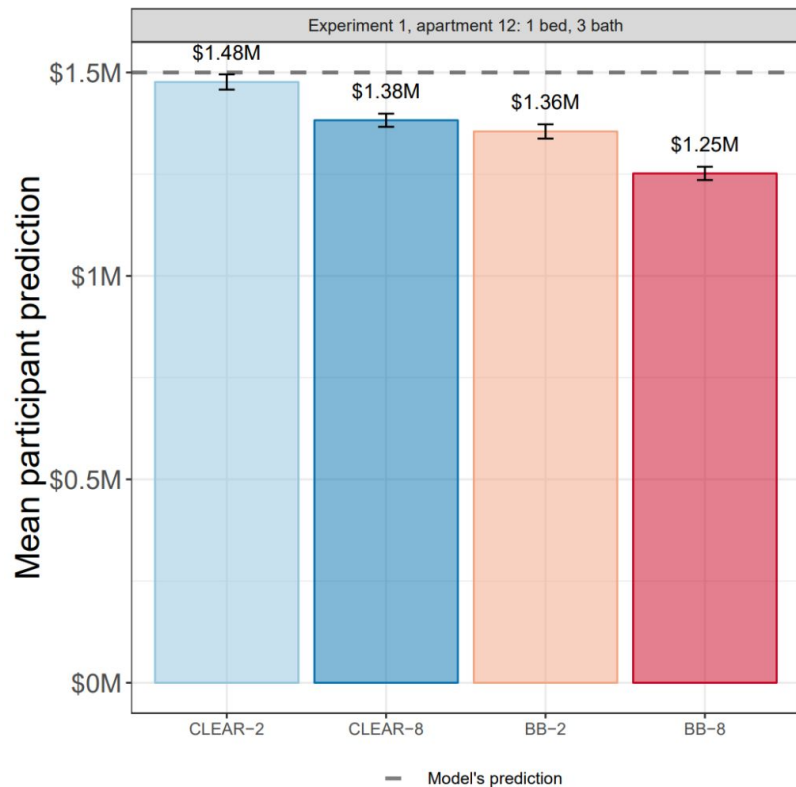


How confident are you that you got it right?



## EXPERIMENT 1: PREDICTING APARTMENT SELLING PRICES (N=1250)







## Experiment 1: New York City prices

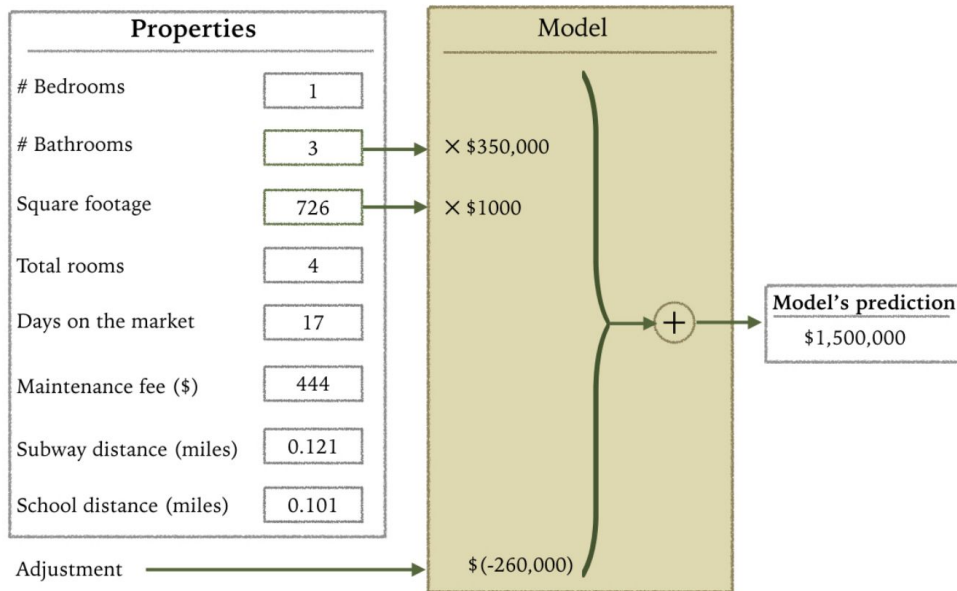


## Experiment 2: Representative U.S. prices



## EXPERIMENT 4: OUTLIER FOCUS AND DETECTION OF MISTAKES (N=800)

**Attention: This apartment has an unusual combination of # Bedrooms and # Bathrooms.**

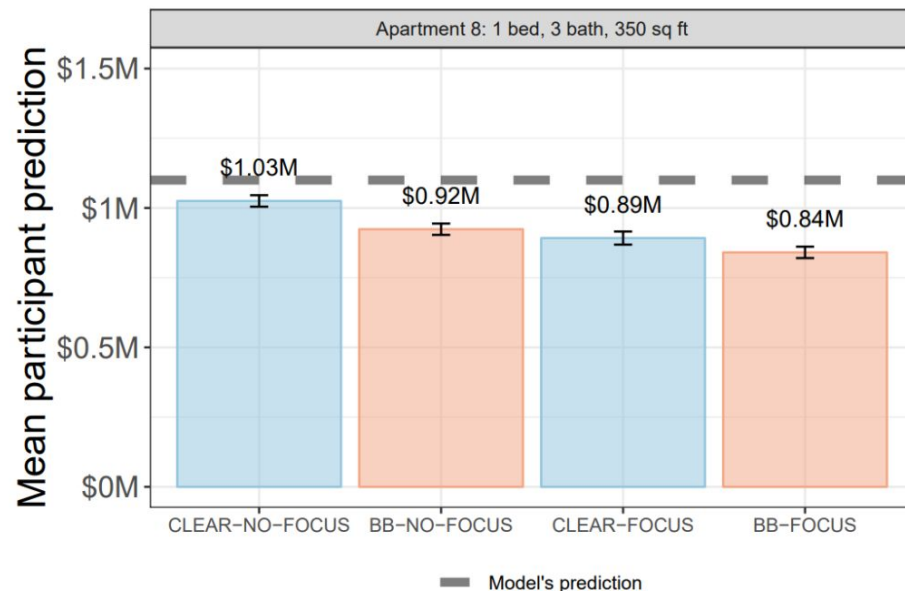
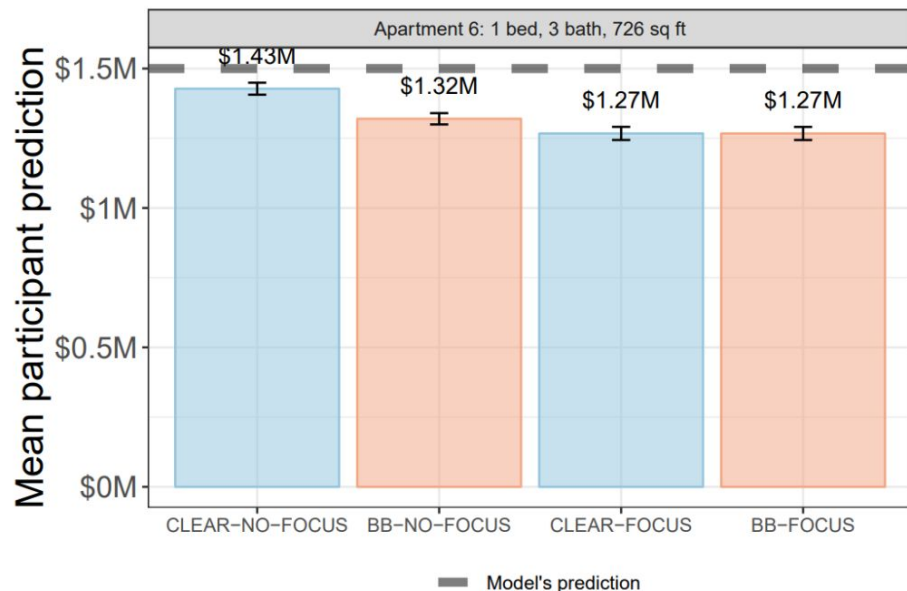


**Two-feature white-box vs  
two-feature black-box**

**Warning vs no warning**

**Please take the unusual configuration of this apartment into consideration when making predictions.**

## EXPERIMENT 4: OUTLIER FOCUS AND DETECTION OF MISTAKES (N=800)



# Hypotheses and outcomes

1. **Outlier focus.** Showing the warning improves performance.
2. **Transparency (clear vs. black box) and no outlier focus.**  
Showing the white-box improves performance.
3. **Transparency (clear vs. black box) and outlier focus.**  
No differences in white-box and black-box.

1. No significant differences (in the users' performance) between a simple white-box and complex black-box

## **2. Information overload**

3. Set goals with respect to interpretability:  
use testing, not intuition

# Theoretical benchmark

Liu et al.

\*Abacus.AI

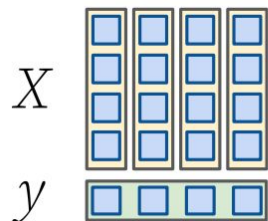
Synthetic Benchmarks for  
Scientific Research in  
Explainable Machine Learning  
(NeurIPS Dataset Track 2021)

---



# XAI-Bench

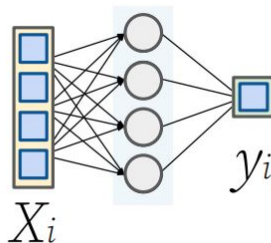
<https://github.com/abacusai/xai-bench>



Data

**Examples:**

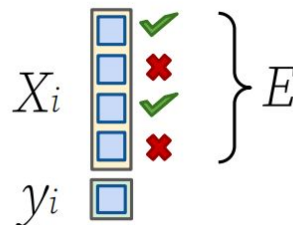
real data,  
synthetic data  
families



Model

**Examples:**

multilayer perceptron,  
decision tree,  
linear regression



Explainer

**Examples:**

SHAP, SHAPR, BF-SHAP,  
MAPLE, LIME, L2X,  
breakDown, Random

$E$	$m$
$E_1$	$m_1$
$E_2$	$m_2$
...	...

Metrics

**Examples:**

GT-shapley, ROAR  
faithfulness,  
monotonicity

# Data

1. **Generate features  $\mathbf{X}$ .** They have known distributions, e.g. multivariate Gaussian, mixture of Gaussians, and multinomial. Moreover, we know the conditional distributions.

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

2. **Generate labels  $\mathbf{Y}$  from  $\mathbf{X}$ .** For example: piecewise linear, additive targets, nonlinear cosine, exponent. We can scale regression tasks to have the target of mean 0/sd 1.

# Metrics

$$\text{faithfulness} = \text{Pearson} \left( \left| \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}(\mathbf{x}_{F \setminus i})} [f(\mathbf{x}')] - f(\mathbf{x}) \right|_{1 \leq i \leq D}, [w_i]_{1 \leq i \leq D} \right)$$

Comparing which feature would have the most impact on the model output when individually changed.

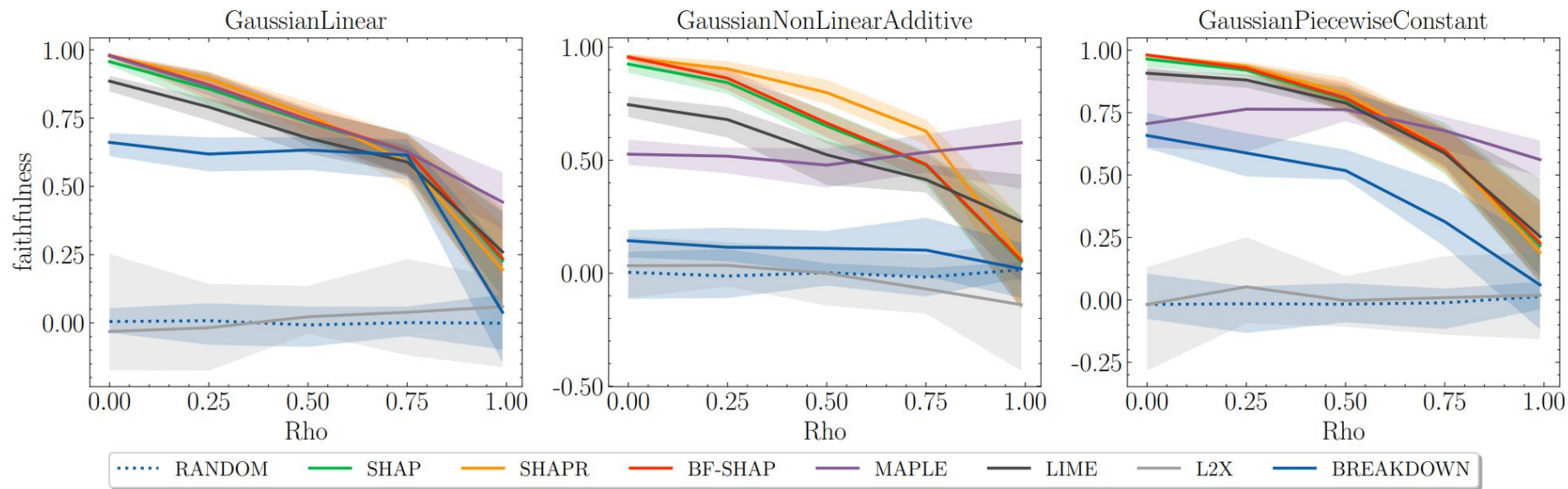
$$\text{monotonicity} = \frac{1}{D-1} \sum_{i=0}^{D-2} \mathbb{I}_{|\delta_i^+| \leq |\delta_{i+1}^+|},$$

$$\text{where } \delta_i^+ = \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}(\mathbf{x}_{S+(w, i+1)})} [f(\mathbf{x}')] - \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}(\mathbf{x}_{S+(w, i)})} [f(\mathbf{x}')]$$

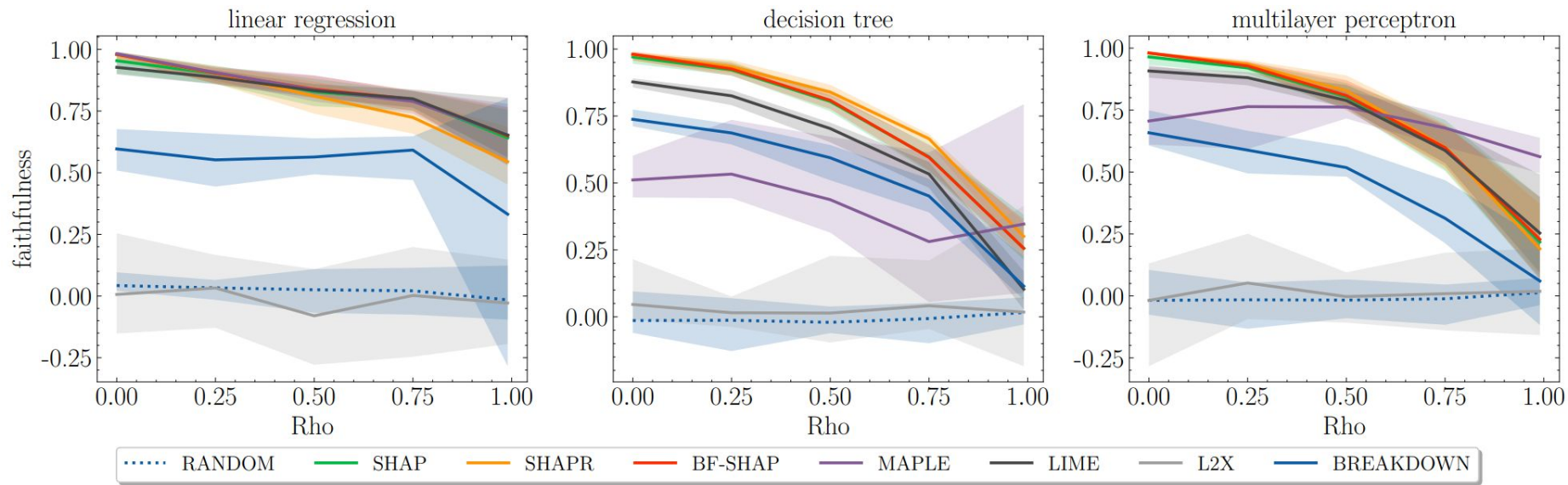
Fraction of indices  $i$  such that the marginal improvement for feature  $i$  is greater than the marginal improvement for feature  $i+1$ .

# Metrics

1. **Remove-and-retrain (ROAR).** Hooker et al. A Benchmark for Interpretability Methods in Deep Neural Networks (NeurIPS 2019)
  - The model is retrained using a new dataset with the features removed.
2. **GT-Shapley.**
  - Pearson between an explanation and ground-truth Shapley values
3. **Infidelity.** Yeh et al. On the (in) fidelity and sensitivity of explanations (NeurIPS 2019)
  - Fidelity but with noisy baseline conditional expectation



**Rho: scale of feature dependence**



**Rho: scale of feature dependence**



# We can approximate the known datasets!

Measure similarities between data distribution  
with Jensen-Shannon Divergence

Table 2: Explainer performance on the simulated wine dataset across metrics. All performance numbers are from explainers for a decision tree.

	RANDOM	SHAP	SHAPR	LIME	MAPLE	L2X	BREAKDOWN
faithfulness ( $\uparrow$ )	$-0.007 \pm 0.005$	<b>0.534</b> $\pm 0.045$	$0.528 \pm 0.032$	$0.368 \pm 0.031$	$0.034 \pm 0.033$	$-0.030 \pm 0.018$	$-0.042 \pm 0.011$
monotonicity ( $\uparrow$ )	$0.529 \pm 0.008$	$0.549 \pm 0.009$	<b>0.551</b> $\pm 0.009$	$0.547 \pm 0.007$	$0.520 \pm 0.014$	$0.522 \pm 0.005$	$0.493 \pm 0.014$
ROAR ( $\uparrow$ )	$0.698 \pm 0.031$	$0.780 \pm 0.016$	$0.549 \pm 0.031$	$0.738 \pm 0.026$	<b>0.818</b> $\pm 0.022$	$0.664 \pm 0.02$	$0.625 \pm 0.002$
GT-Shapley ( $\uparrow$ )	$0.004 \pm 0.013$	$0.825 \pm 0.006$	<b>0.945</b> $\pm 0.002$	$0.745 \pm 0.015$	$0.685 \pm 0.008$	$-0.108 \pm 0.029$	$-0.064 \pm 0.02$
infidelity ( $\downarrow$ )	$0.353 \pm 0.174$	$0.234 \pm 0.124$	<b>0.212</b> $\pm 0.146$	$0.234 \pm 0.126$	$0.234 \pm 0.132$	$0.285 \pm 0.115$	$0.365 \pm 0.133$

Table 3: Explainer performance on the simulated forest fires dataset across metrics. All performance numbers are from explainers for a decision tree.

	RANDOM	SHAP	LIME	MAPLE	L2X	BREAKDOWN
faithfulness ( $\uparrow$ )	$0.022_{\pm 0.034}$	<b>0.571</b> $_{\pm 0.023}$	$0.449_{\pm 0.007}$	$0.080_{\pm 0.056}$	$0.001_{\pm 0.008}$	$0.158_{\pm 0.032}$
monotonicity ( $\uparrow$ )	$0.537_{\pm 0.02}$	$0.591_{\pm 0.007}$	<b>0.598</b> $_{\pm 0.002}$	$0.561_{\pm 0.002}$	$0.527_{\pm 0.01}$	$0.575_{\pm 0.012}$
ROAR ( $\uparrow$ )	$0.575_{\pm 0.002}$	$0.615_{\pm 0.011}$	$0.616_{\pm 0.008}$	<b>0.696</b> $_{\pm 0.024}$	$0.534_{\pm 0.018}$	$0.604_{\pm 0.019}$
GT-Shapley ( $\uparrow$ )	$0.012_{\pm 0.06}$	<b>0.870</b> $_{\pm 0.005}$	$0.779_{\pm 0.027}$	$0.804_{\pm 0.011}$	$0.031_{\pm 0.12}$	$0.105_{\pm 0.013}$
infidelity ( $\downarrow$ )	$0.207_{\pm 0.125}$	<b>0.075</b> $_{\pm 0.074}$	$0.077_{\pm 0.075}$	$0.077_{\pm 0.079}$	$0.091_{\pm 0.07}$	$0.117_{\pm 0.076}$

We should always be able to create artificial data that allows evaluating our algorithms.

# Practical experiment

Zhou & Ribeiro et al.

\*MIT

Do Feature Attribution  
Methods Correctly Attribute  
Features?  
(AAAI 2022)

---

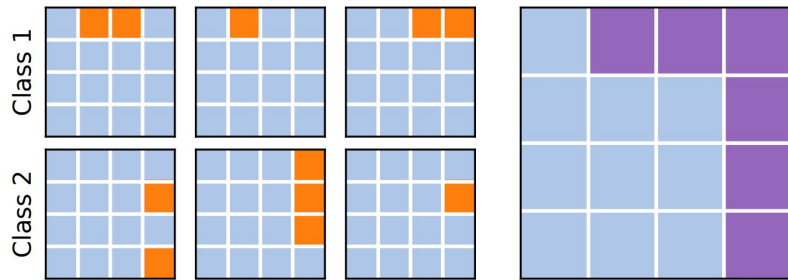


Figure 1: The intuition behind our feature attribution ground truth: if we know that for every input, only specific features (orange) are informative to the label, then across the dataset, a high-performing model has to focus on them and not get “distracted” by other irrelevant features. Thus, feature attributions should highlight the union *union* of these features (purple), and any attribution outside this area is misleading.

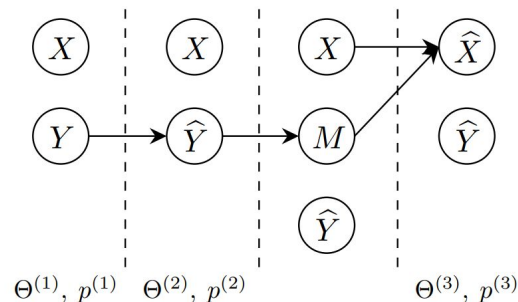
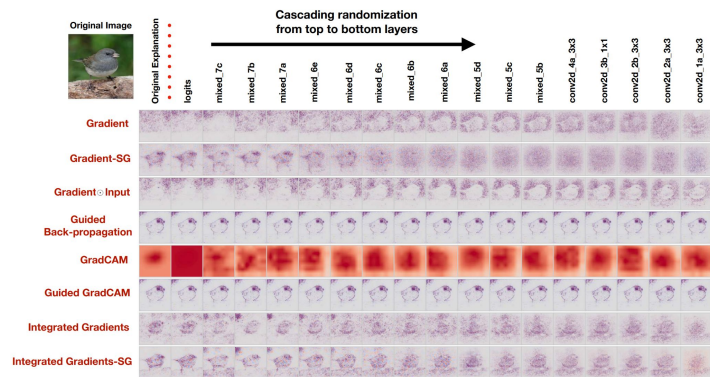


Figure 2: The graphical model for our dataset modification.

We can evaluate explanations against the crafted ground-truth.



# Recap

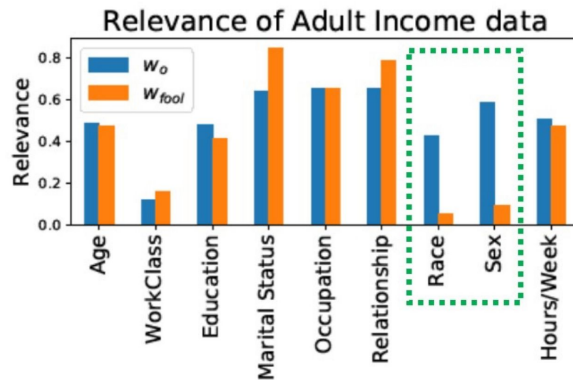
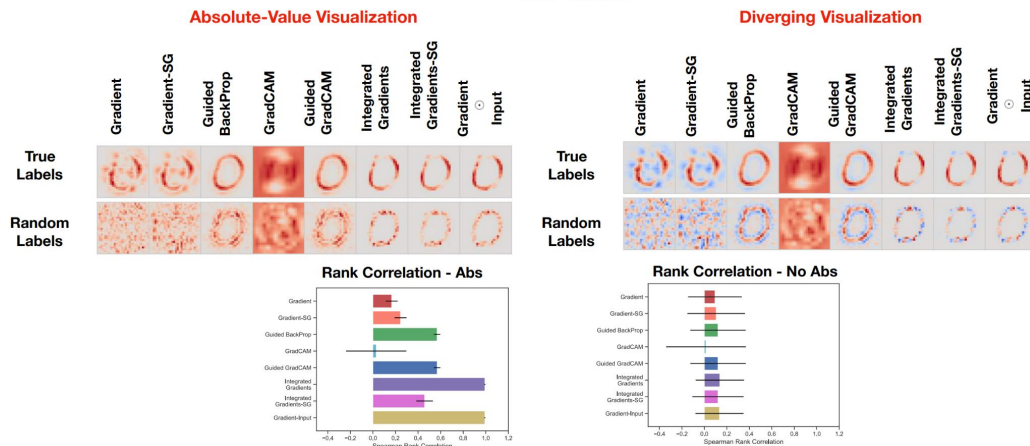


Adebayo et al. Sanity Checks for Saliency Maps (NeurIPS 2018)



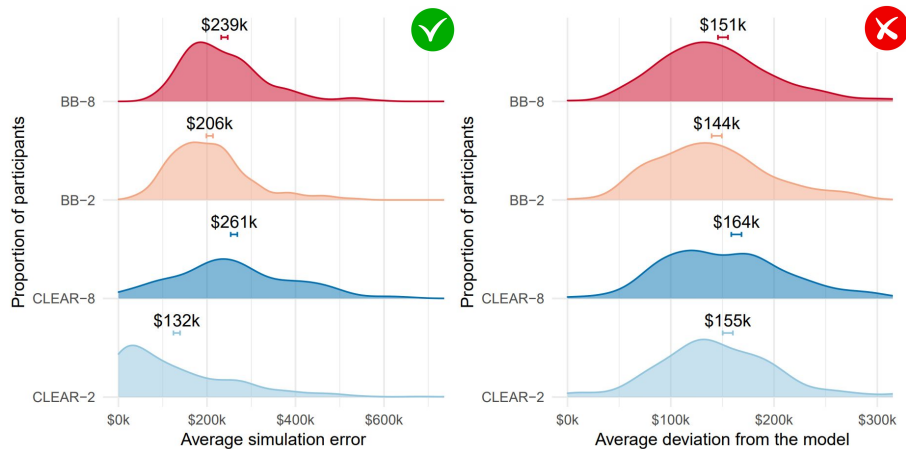
Dombrowski et al. Explanations can be manipulated and geometry is to blame (NeurIPS 2019)

CNN - MNIST

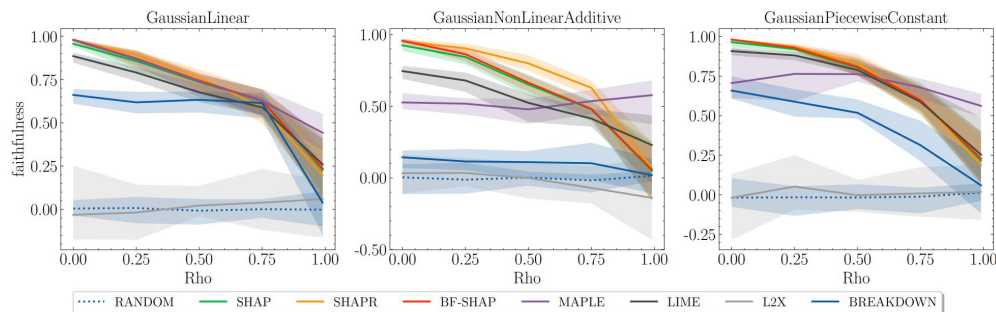


Heo et al. Fooling Neural Network Interpretations via Adversarial Model Manipulation (NeurIPS 2019)

# Recap



Poursabzi-Sangdeh et al. Manipulating and Measuring Model Interpretability (CHI 2021)



Liu et al. Synthetic Benchmarks for Scientific Research in Explainable Machine Learning (NeurIPS Dataset Track 2021)

Zhou et al. Do Feature Attribution Methods Correctly Attribute Features? (AAAI 2022)