

\mathcal{M}^4 : A Unified XAI Benchmark for Faithfulness Evaluation of Feature Attribution Methods across Metrics, Modalities, and Models

Xuhong Li¹, Mengnan Du², Jiamin Chen¹,
Yekun Chai¹, Himabindu Lakkaraju³, Haoyi Xiong¹

¹ Baidu Inc. ² New Jersey Institute of Technology
³ Harvard University

NeurIPS 2023 Datasets and Benchmarks



Faithfulness of XAI Algorithms

- Feature attributions.
- Faithfulness?
- Recent benchmarks.



Faithfulness of XAI Algorithms

- Feature attributions.
- Faithfulness?
- Recent benchmarks.

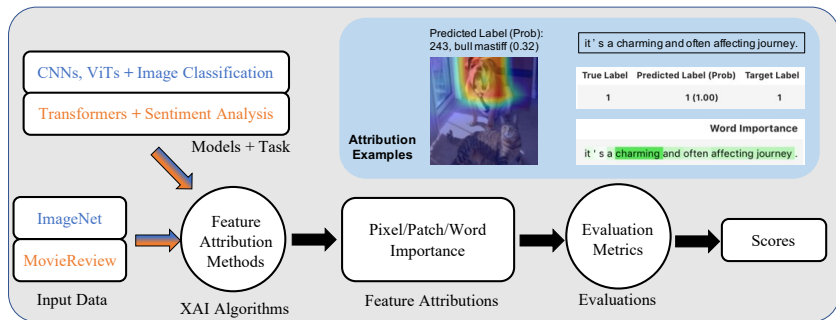


Faithfulness of XAI Algorithms

- Feature attributions.
- Faithfulness?
- Recent benchmarks.



- A unified XAI benchmark for feature attribution Methods across Modalities, Models, and Metrics.



Tasks, Datasets and Models:

- Image classification
 - 5,000 images from ImageNet validation set
 - VGG, ResNets, Mobilenet-V3, ViTs, MAE-ViTs
- Sentiment analysis
 - Movie Review dataset
 - BERTs, DistilBERT, ERNIE, RoBERTa

Feature Attribution Methods:

- Model-agnostic: LIME
- Gradient-based: Integrated Gradient, SmoothGrad, GradCAM
- Transformer-specific: Generic Attribution, Head-wise/Token-wise Bidirectional Transformer Attributions



Metrics and Taxonomy

- No Ground Truth

- **Most Relevant First:** $\text{MoRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{MoRF}}^{(0)}) - f(x_{\text{MoRF}}^{(k)}))$
- **Least Relevant First:** $\text{LeRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{LeRF}}^{(0)}) - f(x_{\text{LeRF}}^{(k)}))$
- **Area Between Perturbation Curves**
- **INFidelity:** $\text{INFD}(x) = E_{l \sim \mu_l} (l^T A(x, f) - (f(x) - f(x - l))^2)$

- Pseudo Ground Truth

- *M* a set of well-trained models:
 $\text{PScore}(x) = \cos(\frac{1}{|M|} \sum_{g \in M} A(x, g), A(x, l))$

- Synthetic Ground Truth

- *metric* \in {AP, AUO, ROC} and *Syn*(*x*) the synthetic ground truth of explanation:
 $\text{SynScore}_{\text{metric}}(x) = \text{metric}(\text{Syn}(x), A(x, f))$

- MoRF + ABPC + PScore – INFD + SynScore



Metrics and Taxonomy

- No Ground Truth

- **Most Relevant First:** $\text{MoRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{MoRF}}^{(0)}) - f(x_{\text{MoRF}}^{(k)}))$
- **Least Relevant First:** $\text{LeRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{LeRF}}^{(0)}) - f(x_{\text{LeRF}}^{(k)}))$
- **Area Between Perturbation Curves**
- **INFiD**elity: $\text{INFD}(x) = E_{l \sim \mu_l} (l^T A(x, f) - (f(x) - f(x - l))^2)$

- Pseudo Ground Truth

- M a set of well-trained models:
PScore $(x) = \cos(\frac{1}{|M|} \sum_{g \in M} A(x, g), A(x, f))$

- Synthetic Ground Truth

- $\text{metric} \in \{\text{AP}, \text{AUC}, \text{ROC}\}$ and $\text{Syn}(x)$ the synthetic ground truth of explanation:
 $\text{SynScore}_{\text{metric}}(x) = \text{metric}(\text{Syn}(x), A(x, f))$

- MoRF + ABPC + PScore – INFD + SynScore



Metrics and Taxonomy

- No Ground Truth
 - **Most Relevant First**: $\text{MoRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{MoRF}}^{(0)}) - f(x_{\text{MoRF}}^{(k)}))$
 - **Least Relevant First**: $\text{LeRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{LeRF}}^{(0)}) - f(x_{\text{LeRF}}^{(k)}))$
 - **Area Between Perturbation Curves**
 - **INFidelity**: $\text{INFD}(x) = E_{l \sim \mu_l} (l^T A(x, f) - (f(x) - f(x - l))^2)$
- Pseudo Ground Truth
 - M a set of well-trained models:
 $\text{PScore}(x) = \cos(\frac{1}{|M|} \sum_{g \in M} A(x, g), A(x, f))$
- Synthetic Ground Truth
 - $metric \in \{\text{AP}, \text{AUC-ROC}\}$ and $\text{Syn}(x)$ the synthetic ground truth of explanation:
 $\text{SynScore}_{metric}(x) = metric(\text{Syn}(x), A(x, f))$
- MoRF + ABPC + PScore – INFD + SynScore

Metrics and Taxonomy

- No Ground Truth
 - **Most Relevant First**: $\text{MoRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{MoRF}}^{(0)}) - f(x_{\text{MoRF}}^{(k)}))$
 - **Least Relevant First**: $\text{LeRF}(x) = \frac{1}{L+1} \sum_{k=0}^L (f(x_{\text{LeRF}}^{(0)}) - f(x_{\text{LeRF}}^{(k)}))$
 - **Area Between Perturbation Curves**
 - **INFiD**elity: $\text{INFD}(x) = E_{l \sim \mu_l} (l^T A(x, f) - (f(x) - f(x - l))^2)$
- Pseudo Ground Truth
 - M a set of well-trained models:
 $\text{PScore}(x) = \cos(\frac{1}{|M|} \sum_{g \in M} A(x, g), A(x, f))$
- Synthetic Ground Truth
 - $metric \in \{\text{AP}, \text{AUC-ROC}\}$ and $\text{Syn}(x)$ the synthetic ground truth of explanation:
 $\text{SynScore}_{metric}(x) = metric(\text{Syn}(x), A(x, f))$
- MoRF + ABPC + PScore – INFD + SynScore

Benchmark \mathcal{M}^4

```
1 import interpretdl as it
2
3 # Load a pretrained model from PaddlePaddle model zoo.
4 from paddle.vision.models import resnet50
5 model = resnet50(pretrained=True)
6
7 # Available feature attribution methods include but are not limited to
8   SG, IG, LIME, BT, GA and etc. 'interpret' is the universal api.
9 algo = it.SmoothGradInterpreter(model, device="gpu:0")
10 expl_result = algo.interpret("test.jpg")
11
12 # Available faithfulness evaluation metrics include but are not
13   limited to MoRF, ABPC, INFD and etc. 'evaluate' is the universal
14   api. Note that some evaluators do not require the model.
15 evaluator = it.Infidelity(model)
16 eval_result = evaluator.evaluate("test.jpg", expl_result)
```

The modular design of \mathcal{M}^4 :

- compatible with PaddlePaddle, Pytorch, Hugging Face, etc.
- easy with new methods, models, and metrics.



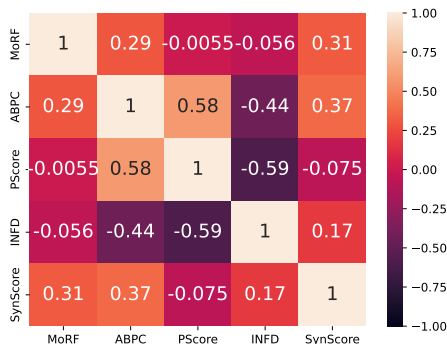
- ▶ **Whether There are Two Metrics that are Correlated ?**
- ▶ Which Explanation Algorithm Demonstrates the Best Faithfulness ?
- ▶ Which Model is the Most (In)sensitive to Explanation Algorithms ?

- ▶ **Whether There are Two Metrics that are Correlated ?**
- ▶ **Which Explanation Algorithm Demonstrates the Best Faithfulness ?**
- ▶ Which Model is the Most (In)sensitive to Explanation Algorithms ?

- ▶ **Whether There are Two Metrics that are Correlated ?**
- ▶ **Which Explanation Algorithm Demonstrates the Best Faithfulness ?**
- ▶ **Which Model is the Most (In)sensitive to Explanation Algorithms ?**

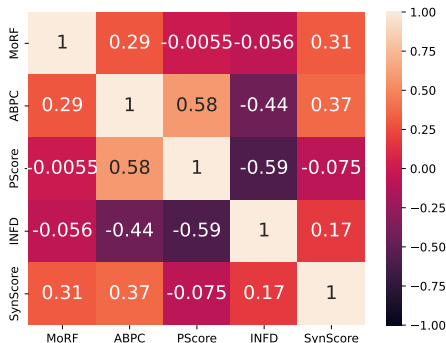
Q1: Correlation between Metrics

- ABPC, PScore, and INFD are strongly correlated.
- MoRF and SynScore evaluate faithfulness from other perspectives.

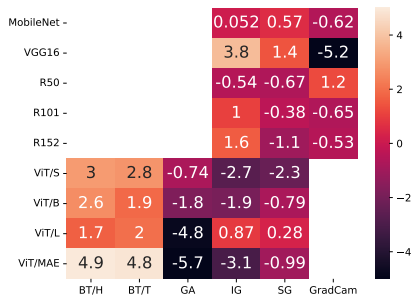


Q1: Correlation between Metrics

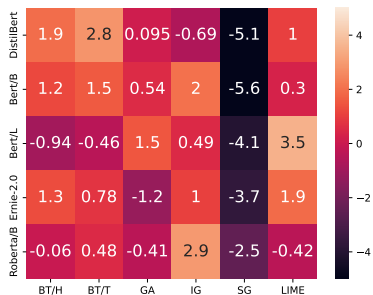
- ABPC, PScore, and INFD are strongly correlated.
- MoRF and SynScore evaluate faithfulness from other perspectives.



Q2: Best Faithfulness

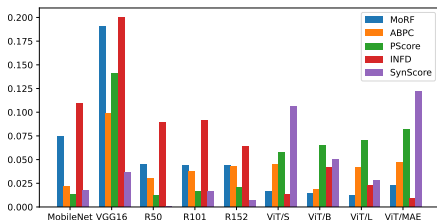
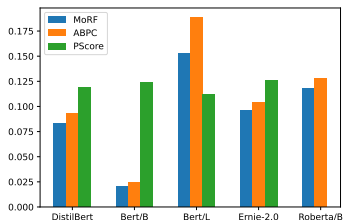


(a) Modality of images.



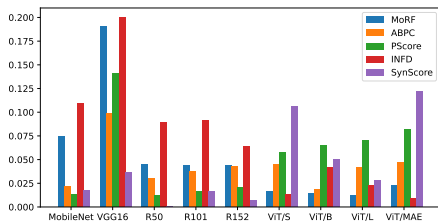
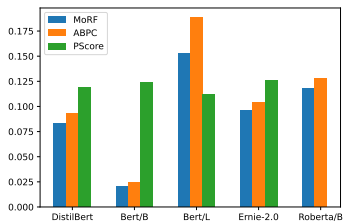
(b) Modality of texts.

Q3: Sensitivity of Models to Explanation Methods



- The sensitivity to different methods can roughly reflect the difficulty of explaining the model.
- Some methods may work especially well for certain models but not as well for others.

Q3: Sensitivity of Models to Explanation Methods



- The sensitivity to different methods can roughly reflect the difficulty of explaining the model.
- Some methods may work especially well for certain models but not as well for others.

The demos, implementations,
and additional information of \mathcal{M}^4
are publicly available at

InterpretDL

<https://github.com/PaddlePaddle/InterpretDL>.

