Causal Intersectionality and Dual Form of Gradient Descent for Multimodal Analysis : a Case Study on Hateful Memes

> <u>Yosuke Miyanishi</u>, Nguyen Le Minh Japan Advanced Institute of Science and Technology (JAIST) LREC-COLING 2024



Table of Contents

Introduction	Our Motivation and Approach / Causality / Explainable AI (XAI) / Causality and XAI	
Related Work	Vision and Language (VL) / Causality in VL / Causal Intersectionality / Attention Attribution Score / Large Language Models (LLMs)	
Experiment I	Metric for Causal Effect / Metric for Model Inner Workings / Experimental Settings / Results (miATE) / Results (MIDAS) / Results (MIDAS vs miATE)	
Experiment II	Objective / Task Design / Experimental Settings / Results (miATE) / Results (MIDAS) / Results (MIDAS vs miATE)	
Conclusion	Conclusion	

Introduction

 $\mathbf{0}$

Our Motivation and Approach / Causality / Explainable AI (XAI) / Causality and XAI

Our Motivation and Approach

Key Question

Is current machine learning framework sufficient for quantifying the model <u>understanding</u> of a task?

What We Do / What We Don't

- We <u>do not</u> create new models or datasets
- We propose the formal definition of a task
- We propose an analytical framework to quantify the relevance between the task's causal structure and inner workings of the SOTA models



- Causality describes the structure of the real world
- Two lines of causality science for graphical repersentation and quantification

Graphical Representation^{*1}





- Causality describes the structure of the real world
- Two lines of causality science for graphical repersentation and quantification

Graphical Representation^{*1}



- Causality describes the structure of the real world
- Two lines of causality science for graphical repersentation and quantification



*1: Pearl, 2000 (10.1017/S0266466603004109) *2: Rubin, 2008 (10.1214/08-AOAS187)

Causality

- Causality describes the structure of the real world
- Two lines of causality science for graphical repersentation and quantification



*1: Pearl, 2000 (10.1017/S0266466603004109) *2: Rubin, 2008 (10.1214/08-AOAS187)



- Causality describes the structure of the real world
- Two lines of causality science for graphical repersentation and quantification



How to relate the model inner working to the real-world causal structure?

*1: Pearl, 2000 (10.1017/S0266466603004109) *2: Rubin, 2008 (10.1214/08-AOAS187)

Explainable AI (XAI)

- XAI describes how the model reacts to the world
- Centric to XAI is gradient-based method

Example in Computer Vision^{*1}



How can XAI method describe the understanding of the causality?

*1: Selvaraju et al., 2017 (10.1109/ICCV.2017.74)

Related Work

02

Vision and Language (VL) / Causality in VL / Causal Intersectionality / Attention Attribution Score / Large Language Models (LLMs)

• Unique interplay of two modalities (text and image)

Visual Question Answering^{*1} Who is wearing glasses?



*1: Goyal et al., 2017 (10.1109/CVPR.2017.670)

- Unique interplay of two modalities (text and image)
- Strong interplay in hateful memes



Benign / Hateful Memes^{*2}

vou smell todav

Love the way

*1: Goyal et al., 2017 (10.1109/CVPR.2017.670) *2: Kiela et al., 2020 (10.48550/arXiv.2005.04790)

- Unique interplay of two modalities (text and image)
- Strong interplay in hateful memes



*1: Goyal et al., 2017 (10.1109/CVPR.2017.670) *2: Kiela et al., 2020 (10.48550/arXiv.2005.04790)

- Unique interplay of two modalities (text and image)
- Strong interplay in hateful memes
 - Visual Question Answering^{*1}

Who is wearing glasses?



 Benign / Hateful Memes*2

 Benign / Love the way

 Journal of the way

We focus on the study of hateful memes

*1: Goyal et al., 2017 (10.1109/CVPR.2017.670) *2: Kiela et al., 2020 (10.48550/arXiv.2005.04790)



• Causality <u>of</u> / <u>to</u> the model has been tested for Visual Question Answering (VQA)

• Causality <u>of</u> / <u>to</u> the model has been tested for Visual Question Answering (VQA)

Bi-Phasic Cognitive VQA Model^{*1}



Causality <u>of</u> / <u>to</u> the model has been tested for Visual Question Answering (VQA)



*1: Nguyen and Okazaki, 2023 (10.18653/v1/2023.emnlp-main.573) *2: Niu et al., 2021 (10.1109/CVPR46437.2021.01251)

• Causality <u>of</u> / <u>to</u> the model has been tested for Visual Question Answering (VQA)



*1: Nguyen and Okazaki, 2023 (10.18653/v1/2023.emnlp-main.573) *2: Niu et al., 2021 (10.1109/CVPR46437.2021.01251)

• Causality <u>of</u> / <u>to</u> the model has been tested for Visual Question Answering (VQA)



*1: Nguyen and Okazaki, 2023 (10.18653/v1/2023.emnlp-main.573) *2: Niu et al., 2021 (10.1109/CVPR46437.2021.01251)

- Interplay of two social categories (e.g. colour and gender)
- Interpreted as a form of indirect causal effect

gender, especially when this may result in additional

disadvantage or discrimination

Definition*1Intersectionality as Indirect Effect*2the network of connections
between social categories
such as race, class, and $\theta_{Z_1,Z_2} \neq \theta_{Z_1} + \theta_{Z_2}$

- Interplay of two social categories (e.g. colour and gender)
- Interpreted as a form of indirect causal effect



 Z_1

 Z_2

- Interplay of two social categories (e.g. colour and gender)
- Interpreted as a form of indirect causal effect



 Z_1

 Z_2

disadvantage or discrimination

Interplay of two social categories (e.g. colour and gender)
Interpreted as a form of indirect causal effect



We expand to analyze multimodal indirect effect

Gradient-Based Method in NLP

• Attention Attribution Score^{*1}: attention importance to the model prediction



*1: Hao et al., 2021 (10.1609/aaai.v35i14.17533) *2: Sundararajan et al., 2017 (10.5555/3305890.3306024)

Gradient-Based Method in NLP

- Attention Attribution Score^{*1}: attention importance to the model prediction
- Hee et al.^{*3} showed modality attribution



*1: Hao et al., 2021 (10.1609/aaai.v35i14.17533) *2: Sundararajan et al., 2017 (10.5555/3305890.3306024) *3: Hee et al., 2022 (10.1145/3485447.3512260) *4: Li et al., 2020 (10.18653/v1/2020.acl-main.469)

Gradient-Based Method in NLP

- Attention Attribution Score^{*1}: attention importance to the model prediction
- Hee et al.^{*3} showed modality attribution



*1: Hao et al., 2021 (10.1609/aaai.v35i14.17533) *2: Sundararajan et al., 2017 (10.5555/3305890.3306024) *3: Hee et al., 2022 (10.1145/3485447.3512260) *4: Li et al., 2020 (10.18653/v1/2020.acl-main.469)

Large Language Models (LLMs)

Notable performance on gradient-free in-context learning (ICL) setting



*1: Brown et al., 2020 (10.5555/3495724.3495883)

Large Language Models (LLMs)

- Notable performance on gradient-free in-context learning (ICL) setting
- Meta-gradient/optimization: attention weights as the second form of gradient



*1: Brown et al., 2020 (10.5555/3495724.3495883) *2: Dai et al., 2023 (10.18653/v1/2023.findings-acl.247)

Large Language Models (LLMs)

- Notable performance on gradient-free in-context learning (ICL) setting
- Meta-gradient/optimization: attention weights as the second form of gradient



We show how ICL and meta-gradient contributes to causal effect

*1: Brown et al., 2020 (10.5555/3495724.3495883) *2: Dai et al., 2023 (10.18653/v1/2023.findings-acl.247)

Experiment I

 $\mathbf{03}$

Metric for Causal Effect / Metric for Model Inner Workings / Experimental Settings / Results (miATE)

• Extended causal intersectionality for multimodal analysis

Causal Intersectionality / ATE







• Extended causal intersectionality for multimodal analysis



• Extended causal intersectionality for multimodal analysis



• Extended causal intersectionality for multimodal analysis



Metric for Model Inner Workings

Modality Interaction Disentangled Attribution Score (MIDAS)

Attr(A)Interaction Type it \in {*text2text, image2image, cross modal*} Attr and = A $\partial \theta(\alpha A) d\alpha$ MIDAS^{it} MIDAS $= Attr^{it}(A_{T,I}) - (Attr^{it}(A_T) + Attr^{it}(A_I))$

Experimental Settings



*1: Muennighoff, 2020 (arXiv:2012.07788) *2: Sundararajan et al., 2017 (10.5555/3305890.3306024) *3: Li et al., 2020 (10.1007/978-3-030-58577-8_8)

*4: Guolin et al., 2017 (10.5555/3294996.3295074) *5: Akiba et al., 2019 (10.1145/3292500.3330701)

Results (miATE)

- Higher miATE in text-oriented task
- Most remarkable in VisualBERT



More sensitive to the difference of text?

• In contrast to Attr, captured the attention to the difference of modalities



• In contrast to Attr, captured the attention to the difference of modalities



Oscar



• In contrast to Attr, captured the attention to the difference of modalities



Oscar

• In contrast to Attr, captured the attention to the difference of modalities



Oscar

• In contrast to Attr, captured the attention to the difference of modalities

• VisualBERT is biased toward textual information, enhanced by encoder



VisualBERT



In contrast to Attr, captured the attention to the difference of modalities

VisualBERT is biased toward textual information, enhanced by encoder

Org. Image 1.20E-04 1.00E-04 8.00E-05 -6.00E-05 -4.00E-05 -2.00E-05 0.00E+00 ■VisualBERT VisualBERT (Text-Only)

VisualBERT



Org. Text

• In contrast to <u>Attr, captured the attention to the difference of modalities</u>

• VisualBERT is biased toward textual information, enhanced by encoder

Org. Image

VisualBERT



• In contrast to Attr, captured the attention to the difference of modalities

• VisualBERT is biased toward textual information, enhanced by encoder

Org. Image

VisualBERT



Results (MIDAS vs miATE)

Formal Relationship of miATE and MIDAS

 $\begin{array}{c|c} miATE \\ = \theta_{T,I} \\ - (\theta_{S_T} \\ + \theta_I) \end{array} & \begin{array}{c} Attr(A) \approx A \ast G(A) \text{ where } G(A) = \frac{\partial \theta(A)}{\partial A} \\ MIDAS \approx A_{T,I} \ast G(A_{T,I}) - (A_T \ast G(A_T) + A_I \ast G(A_I)) \\ \sum_n MIDAS \approx \mathbb{E}[A_{T,I}] - (\mathbb{E}[A_T] + \mathbb{E}[A_I]) \end{array}$

Results (MIDAS vs miATE)

Empirical Relationship of miATE and MIDAS

		Oscar	VisualBERT
AUC(%)		74.3 ± 2.39	94.1 ± 4.20
Feature Importance	text2text	35 ± 25	253 <u>+</u> 66
	image2image	30 ± 23	169 ± 74
	cross-modal	19 <u>±</u> 14	169 ± 74



04

Objective and Task Design / Meta Optimization / Experimental Settings / Results

Objective and Task Design

• Classifier's objective implicitly maximize miATE, while zero-shot LLM not



Objective and Task Design

• Classifier's objective implicitly maximize miATE, while zero-shot LLM not

• Presenting positive / negative examples to align with the classifiers



Meta Optimization

• Multifaceted nature of the task for chat-bot style LLM

Task Type Classification (TTC) / Label Identification (LI)



Meta Optimization

• Multifaceted nature of the task for chat-bot style LLM

Task Type Classification (TTC) / Label Identification (LI)



Meta Optimization

- Multifaceted nature of the task for chat-bot style LLM
- Similar analysis for meta-gradient

Task Type Classification (TTC) / Label Identification (LI)



Meta-Gradient by Interaction Type

SubTask st = TTC Interaction Type it $\in \{t2t, i2i, cross modal\}$ $\theta = \sum_{st} \sum_{it} (A_{it}^{st} + \Delta A_{it}^{st}) * Q$

Experimental Settings

Model

<Hateful Memes Detection>
Llama-2^{*1} with BLIP-2^{*2} caption extraction
<Probing for meta-gradient vs subtask output>
LightGBM^{*3} with Optuna^{*4} hyperparameter search

*1: Touvron et al., 2023 (10.48550/arXiv.2307.09288) *2: Li et al., 2020 (10.5555/3618408.3619222) *3: Guolin et al., 2017 (10.5555/3294996.3295074) *4: Akiba et al., 2019 (10.1145/3292500.3330701)

Results

• ICL impacts comprehension of the challenging task

TTC Performance (Accuracy %)



Results

- ICL impacts comprehension of the challenging task
- No big difference of contribution among interaction type



$$\theta = \sum_{st} \sum_{it} \left(A_{it}^{st} + \Delta A_{it}^{st} \right) * Q$$







Conclusion

Conclusion

Contribution

- Causal intersectionality for multimodal analysis
- Metrics for causal effect and model inner workings
- Relationship between causal effect and model inner workings
- Impact of meta-optimization on the model performance

Future Work

Generalizability for other hateful memes datasets
 Generalizability for other problems (e.g. missing modality, medical diagnosis)

Limitations

Resources for hateful memes datasetLanguage constraint

Thanks!

Do you have any questions?

yosuke.miyanishi@jaist.ac.jp

CREDITS: This presentation template was created by <u>Slidesgo</u>, and includes icons by <u>Flaticon</u>, and infographics & images by **Freepik**