

Adversarial Perturbations and Feature Attribution Consistency in Integrated Gradients vs Attention Rollout

Assignee Research

May 31, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: To what extent do adversarial perturbations in input text degrade the consistency of feature attribution maps generated by Integrated Gradients compared to Attention Rollout. Attribution algorithms are frequently employed to explain the decisions of neural network models. Integrated Gradients (IG) is an influential attribution method due to its strong axiomatic foundation. 17 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Integrated Decision Gradients: Compute Your Attributions Where the Model Makes Its Decision. Research question: To what extent do adversarial perturbations in input text degrade the consistency of feature attribution maps generated by Integrated Gradients compared to Attention Rollout?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.3/10.

3 Results

16 papers retrieved. 17 claims extracted; 4 independently verified. Quality review score: 5.3/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The experiments were performed using the 2012 validation set of ImageNet on NVIDIA A40 GPUs.	×	0.03
The attributions computed using Algorithm 1 are called IDG.	×	0.13
IDG is compared with IG, LIG, GIG, and AGI.	×	0.04
Captum is used for the implementation of IG.	×	0.03
LIG, GIG, and AGI are taken from their respective repositories.	×	0.00
The quality of the computed attributions is evaluated both quantitatively and qualitatively.	×	0.06
Standard perturbation testing measures the importance of the pixels in an attribution via an area under the curve (AUC)	×	0.03
Four tests are used with three insertion butions than IG.	×	0.04
Integrated decision gradients is the solution to the saturation problem.	✓	0.27
Adaptive sampling is used to provide IDG access to gradients of a higher quality than uniform sampling does.	×	0.08
IDG improves in both qualitative and quantitative results compared with IG, LIG, GIG, and AGI.	×	0.03
IDG focuses on integrating gradients from the decision region of the path integral.	✓	0.23
IDG scales each gradient by the derivative of the output logit with respect to the path.	✓	0.24
The scaling factor rewards gradients in the decision region and penalizes gradients from outside the decision region.	×	0.10
IDG provides a principled solution to saturation by satisfying the IG axioms and a new path integral sensitivity axiom.	×	0.14
An adaptive sampling technique is presented to select non-uniform subdivisions for the Riemann approximation of the path	✓	0.15
The non-uniform subdivisions reduce computational errors (and runtime overheads) compared with using uniform subdivision	×	0.08

References

- <http://arxiv.org/abs/1801.04693v1>
- <http://arxiv.org/abs/1911.11746v1>
- <http://arxiv.org/abs/2305.20052v2>