

ExpProof : Operationalizing Explanations for Confidential Models with ZKPs

Chhavi Yadav*\$, Evan Monroe Laufer*#, Dan Boneh#, Kamalika Chaudhuri\$ \$UC San Diego, #Stanford University *egual contribution



Explanations failure mode

- Explanations are intended as a way to increase trust in ML models by making them transparent in societal applications and are often obligated by regulations.
- Many of these use-cases are adversarial in nature where the involved parties have misaligned interests and are incentivized to manipulate explanations to meet their ends.
- Consequently, despite the demand, explanations fail to be operational as a trust-enhancing tool.

Confidentiality aggravates the problem

- Additionally, organizations keep their models confidential due to IP & legal reasons.
- Leads to model swapping : for inference vs. explanation, for different inputs and post-audits. Moreover, there is no guarantee over how the explanation is generated.
- Canonical Solution : Consistency checks. But these have been shown to be hard. Require the customer to collect different (input, explanation) pairs which makes for a lopsided ask.

Our Solution

- Public verification of explanations with cryptography
- We propose *ExpProof* a system that uses Zero-Knowledge Proofs (ZKPs) & Commitments to publicly verify explanations, while maintaining confidentiality.
- We also propose ZKP-efficient versions of the explanation algorithm.



Table 1. ZKP Overhead of BorderLIME and Standard LIME (both G+N variant) for NNs. Overhead for BorderLIME is larger than that for LIME. Results are consistent across all datasets.

LIME

 1.17 ± 10^{-2}

 0.11 ± 10^{-2}

 10.40 ± 0

Fidelity of Standard LIME & BorderLIME variants





Is ExpProof computationally feasible for LIME?



17. end for 18: $x_{border} := x_{border_i}$ such that $i := \arg \min dist_i$

19: Return xborder