

Understanding Data Importance in Machine Learning Attacks: Does Valuable Data Pose Greater Harm?

Rui Wen, Michael Backes, Yang Zhang CISPA



Machine learning models are widely deployed























Quantitatively measure the contribution



$$u_{\mathrm{shap}}(z) \propto rac{1}{N} \sum_{S \subseteq D \setminus \{z\}} rac{1}{\binom{N-1}{|S|}} \left[U_{\mathcal{A}, D_{val}}(S \cup \{z\}) - U_{\mathcal{A}, D_{val}}(S)
ight]$$

High importance data contributes much more











Model

8

Data shows different levels of vulnerability







High vulnerability?





records with *rare* but *crucial* symptoms high-importance







Discrimination

Premiums





Vulnerability increases for high importance samples







High vulnerability!









Low importance samples are harder to learn



It's hard to identify low-importance members



Compare samples with similar importance







Calibrate membership metric by sample importance

$$= \operatorname{OriMem}(x) + k \times \operatorname{Shapley}(x)$$

$$0.25$$

$$0.20$$

$$0.15$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

$$0.00$$

CaliMem(x)



Take away: high-importance data is *more vulnerable* to membership inference attacks

Setting **sample-specific thresholds** based on importance can make attacks even stronger.



Removing the "layer" of outlier points that are most vulnerable to a privacy attack exposes a new layer of previously-safe points to the same attack.



[1] The Privacy Onion Effect: Memorization is Relative





[1] Some Results on Privacy and Machine Unlearning, Matthew Jagielski









More important





More important

Privacy onion effect holds for importance value

















Target sample













Activately modify importance can lead to stronger attack



[1] Truth Serum: Poisoning Machine Learning Models to Reveal Their Secrets



Take away: "privacy onion effect" holds for data importance.

Actively manipulating sample importance can be a potent strategy for developing stronger attacks.









