Secure Systems Group, University of Waterloo⁺, Aalto University[†]

Vasisht Duddu⁺, Sebastian Szyller⁺, N. Asokan^{+,†}

SoK: Unintended Interactions among Machine Learning Defenses and Risks

Motivation

• ML models susceptible to different risks to security, privacy, and fairness

Situating prior work in framework

- Risk increases \rightarrow , decreases \rightarrow
- Defenses designed against specific risks

But may also impact unrelated risks Unintended interactions

No systematic framework to understand them

Unintended interactions

Defenses	Risks
RD1 (Adversarial Training) RD2 (Outlier Removal)	R1 (Evasion) R2 (Poisoning)
RD3 (Watermarking) RD4 (Fingerprinting)	R3 (Unauthorized Ownership)
PD1 (Differential Privacy)	P1 (Membership Inference) P2 (Data Reconstruction) P3 (Attribute Inference) P4 (Distribution Inference)

- Interaction unexplored \rightarrow
- Factors evaluated: empirical $\rightarrow \bullet$, theoretical $\rightarrow \odot$, conjectured \rightarrow



Guideline for conjectures

Defences (< \uparrow or \downarrow >, <f>)</f>		
RD1 (Adversarial Training):		
• D1 \uparrow , $ \mathcal{D}_{tr} $ [170]		
• $D2 \downarrow$, tail length [16], [75]		
• $D4 \uparrow$, priority for learning stable attributes [170]		
• 01 ↑, curvature smoothness [108]		

FD1 (Group Fairness) **FD2 (Explanations)**

F (Discriminatory Behaviour)

Conjectured causes: overfitting, memorization

Framework: Underlying causes

Overfitting: Difference in train and test accuracy

Factors: Trainset size (D1); Model capacity (M1)

Memorization: Difference in model prediction on data record w/ and w/o it in training dataset

Influencing factors:

- **Dataset**: Tail length of distribution (D2); \bullet number of attributes (D3); priority of learning stable attributes (D4)
- **Objective function:** curvature smoothness \bullet

● 02.1 ↑, distinguishability in dat	ta records inside and outside \mathcal{D}_{tr} [152]	
• $03 \uparrow$, distance to boundary for	most \mathcal{D}_{tr} data records [185]	
• M1 \uparrow , model capacity [108]		
DD2 (Outlier Demovel).	Risks (< \uparrow or \downarrow >, <f>)</f>	
KD2 (Outlief Kellioval).	R1 (Evasion):	
• D2 ↑, tail length [175]	• D2 ↑, tail length [96], [182]	
RD3 (Watermarking):	• $01 \downarrow$, curvature smoothness [108]	
• D2 ↑, tail length [102]	• 03 \downarrow , distance of \mathcal{D}_{tr} data records to boundary [171]	
• 02.3 \downarrow , distinguishability in o	$_{0.8}$ R2 (Poisoning): atermarks between f_{θ}	
and f_{A}^{der} , but distinct from inde	●● D2 ↑, tail length [17], [102], [127]	
• M1 ↑, model capacity [3]	• M1 ↑, model capacity [3]	
R3 (Unauthorized Model Ownership):		
	• M1 ↓, model capacity [93], [124]	
	P1 (Membership Inference):	
	• D1 \downarrow , $ \mathcal{D}_{tr} $ [144], [193]	
	 D2 ↑, tail length [25], [26] 	
	• D4 \downarrow , priority for learning stable attributes [109], [164]	
	• 02.1 \uparrow , distinguishability for data records inside and outside \mathcal{D}_{tr} [144]	
	• 03 \downarrow , distance to decision boundary [145]	
	• M1 ↑, model capacity [48], [152]	

Effectiveness of defense correlates with factor Change in factor (<f>) correlates with risk

- Use arrows for <defense, f> and <f, risk>:
- If (\uparrow,\uparrow) or $(\downarrow,\downarrow) \rightarrow \bigcirc$; else (\uparrow,\downarrow) or $(\downarrow,\uparrow) \rightarrow \bigcirc$

(O1); distinguishability of observables across datasets (O2.1), subgroups (O2.2), models (O2.3); distance to decision boundary (O3)

Model: same as M1 \bullet

Conjecture is:

- unanimous if all factors agree, or
- determined by dominant factor (O1, O2, O3)

Non-common factors may affect interaction

Aalto University





ssg-research.github.io/mlsec/interactions