

Vasisht Duddu<sup>+</sup>, Sebastian Szyller<sup>†</sup>, N. Asokan<sup>+,†</sup>

# SoK: Unintended Interactions among Machine Learning Defenses and Risks

## Motivation

- ML models **susceptible** to different risks to **security**, **privacy**, and **fairness**
- Defenses designed against specific **risks**

May increase or decrease unrelated risks  
**Unintended interactions**

- No systematic framework** to understand them

## Unintended interactions

Defenses	Risks
RD1 (Adversarial Training) RD2 (Outlier Removal)	R1 (Evasion) R2 (Poisoning) R3 (Unauthorized Ownership)
RD3 (Watermarking) RD4 (Fingerprinting)	P1 (Membership Inference) P2 (Data Reconstruction) P3 (Attribute Inference) P4 (Distribution Inference)
PD1 (Differential Privacy)	F (Discriminatory Behaviour)
FD1 (Group Fairness) FD2 (Explanations)	

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); number of attributes (D3); priority of learning stable attributes (D4)
- Objective function:** curvature smoothness (O1); distinguishability of observables across datasets (O2.1), subgroups (O2.2), models (O2.3); distance to decision boundary (O3)
- Model:** same as M1

## Situating prior work in framework

- Risk increases → ●, decreases → ●
- Interaction unexplored → ●
- Factors evaluated: empirical → ●, theoretical → ○, conjectured → ○

Defenses	Risks	I	OVFT	D1	D2	D3	D4	O1	O2	O3	Both	M1	References
RD1 (Adversarial Training)	R1 (Evasion) R2 (Poisoning) R3 (Unauthorized Model Ownership) P1 (Membership Inference) P2 (Data Reconstruction) P3 (Attribute Inference) P4 (Distribution Inference) F (Discriminatory Behaviour)	● ● ● ● ● ● ● ●	○ ○ ○ ○ ○ ○ ○ ○	● ○ ○ ○ ○ ○ ○ ○	● ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	[96], [108], [182], [202] [161], [179] [90] ( [101] ) [71], [152] [117], [204] [156] [16], [38], [75], [105]
RD2 (Outlier Removal)	R1 (Evasion) R2 (Poisoning) R3 (Unauthorized Model Ownership) P1 (Membership Inference) P2 (Data Reconstruction) P3 (Attribute Inference) P4 (Distribution Inference) F (Discriminatory Behaviour)	● ● ● ● ● ● ● ●	○ ○ ○ ○ ○ ○ ○ ○	● ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○ ○ ○	[63] [163] [26], [47] [82] [142]

## Guideline for conjectures

Defences (<↑ or ↓>, <f>)	Risks (<↑ or ↓>, <f>)
<b>RD1 (Adversarial Training):</b> <ul style="list-style-type: none"><li>D1 ↑,  D<sub>tr</sub>  [170]</li><li>D2 ↓, tail length [16], [75]</li><li>D4 ↑, priority for learning stable attributes [170]</li><li>O1 ↑, curvature smoothness [108]</li><li>O2 .1 ↑, distinguishability in data records inside and outside D<sub>tr</sub> [152]</li><li>O3 ↑, distance to boundary for most D<sub>tr</sub> data records [185]</li><li>M1 ↑, model capacity [108]</li></ul> <b>RD2 (Outlier Removal):</b> <ul style="list-style-type: none"><li>D2 ↑, tail length [175]</li></ul> <b>RD3 (Watermarking):</b> <ul style="list-style-type: none"><li>D2 ↑, tail length [102]</li><li>O2 .3 ↓, distinguishability in data records inside and outside D<sub>tr</sub>, but distinct from index [102]</li><li>M1 ↑, model capacity [3]</li></ul>	<b>R1 (Evasion):</b> <ul style="list-style-type: none"><li>D2 ↑, tail length [96], [182]</li><li>O1 ↓, curvature smoothness [108]</li><li>O3 ↓, distance of D<sub>tr</sub> data records to boundary [171]</li></ul> <b>R2 (Poisoning):</b> <ul style="list-style-type: none"><li>D2 ↑, tail length [17], [102], [127]</li><li>M1 ↑, model capacity [3]</li></ul> <b>R3 (Unauthorized Model Ownership):</b> <ul style="list-style-type: none"><li>M1 ↓, model capacity [93], [124]</li></ul> <b>P1 (Membership Inference):</b> <ul style="list-style-type: none"><li>D1 ↓,  D<sub>tr</sub>  [144], [193]</li><li>D2 ↑, tail length [25], [26]</li><li>D4 ↓, priority for learning stable attributes [109], [164]</li><li>O2 .1 ↑, distinguishability for data records inside and outside D<sub>tr</sub> [144]</li><li>O3 ↓, distance to decision boundary [145]</li><li>M1 ↑, model capacity [48], [152]</li></ul> <b>P2 (Data Reconstruction):</b>

Effectiveness of defense correlates with factor

Change in factor (<f>) correlates with risk

- ↑: **positive** correlation; ↓: **negative** correlation

Use arrows for <defense, f> and <f, risk>:

- If (↑,↑) or (↓,↓) → ●; else (↑,↓) or (↓,↑) → ●

**Conjecture is:**

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

Non-common factors may affect interaction