

Explaining Vulnerabilities to Adversarial Machine Learning Through Visual Analytics

Yuxin Ma, Tiankai Xie, Jundong Li, Ross Maciejewski

Vulnerabilities in Machine Learning

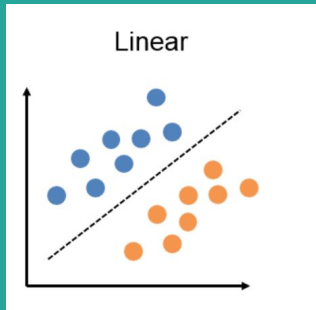


Filtering Spam Emails

Training Stage



Spam/Non spam
Database



Testing Stage



Unlabeled

New
Labels

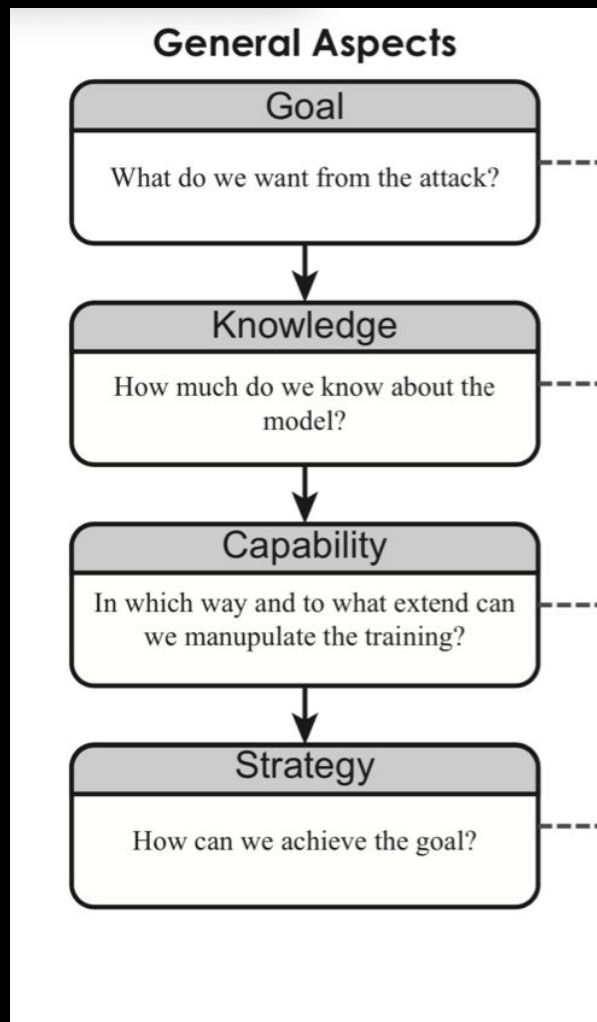


The Contribution of This Paper is:

- A visual analytics framework that supports the examination, creation, and exploration of adversarial machine learning attacks;
- A visual representation of model vulnerability that reveals the impact of adversarial attacks in terms of model performance, instance attributes, feature distributions, and local structures.

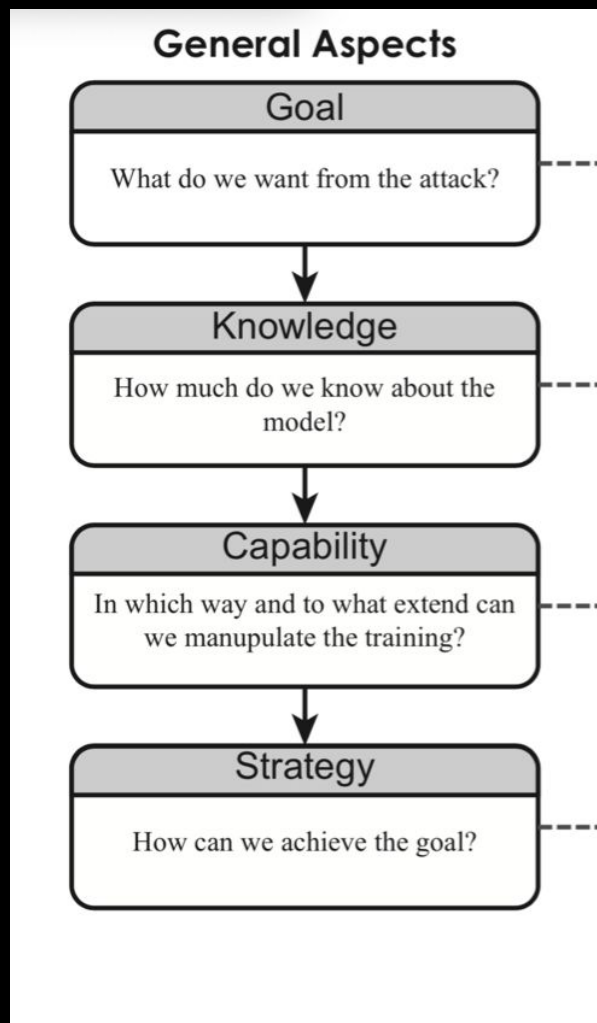
The four main features of an adversary are the adversary's ...

- **Goal**
 - Targeted attack
 - Reliability attack
- Knowledge
- Capability
 - Poisoning attack
 - Evasion attack
- Strategy



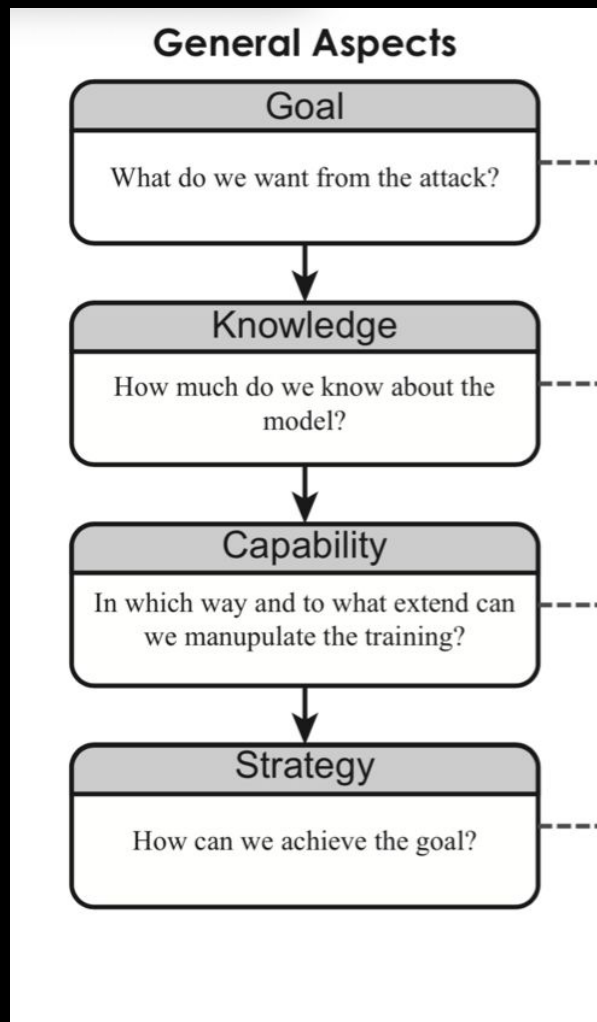
The four main features of an adversary are the adversary's ...

- Goal
 - Targeted attack
 - Reliability attack
- Knowledge
 - Black box
 - White box
- Capability
- Strategy



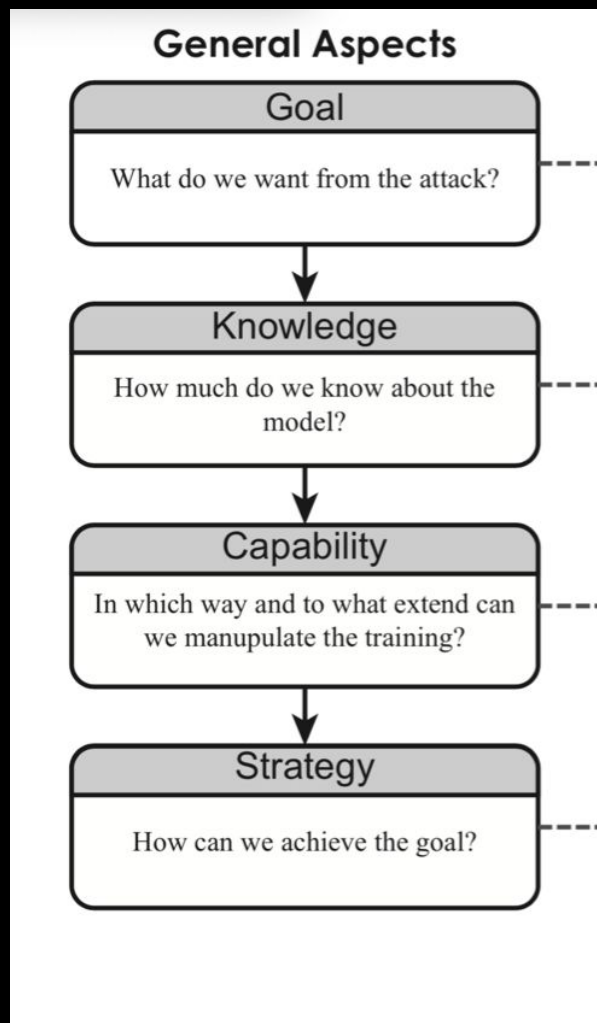
The four main features of an adversary are the adversary's ...

- Goal
 - Targeted attack
 - Reliability attack
- Knowledge
- **Capability**
 - **Poisoning attack**
 - **Evasion attack**
- Strategy



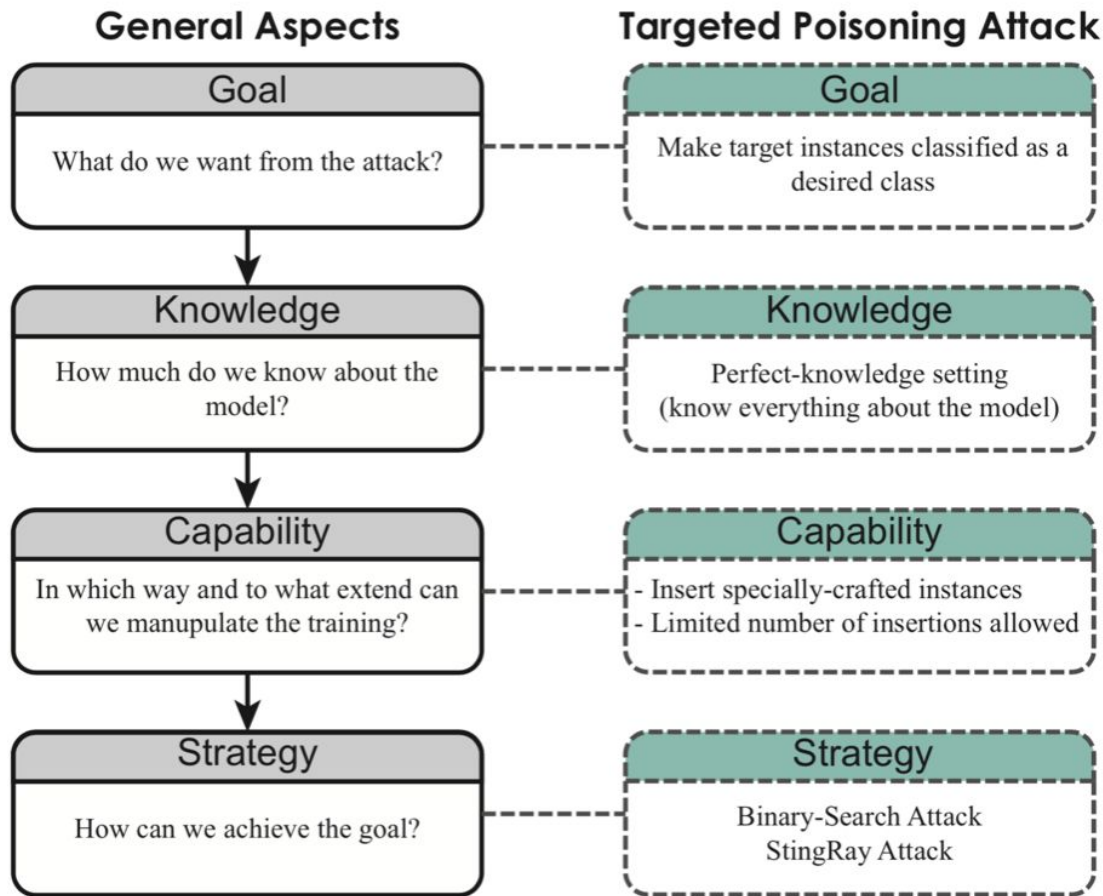
The four main features of an adversary are the adversary's ...

- Goal
 - Targeted attack
 - Reliability attack
- Knowledge
- Capability
 - Poisoning attack
 - Evasion attack
- Strategy



To demonstrate the proposed visual analytics framework, we focus our discussion on:

Targeted Data Poisoning Attack



To demonstrate the

prop

anal

fram

focus

discu

Target

Pois

- **We have summarized two major tasks:**
 - 1. Identify vulnerabilities in the training dataset**
 - 2. Inspect and diagnose what happens when the poisoning instances are inserted.**

How can we achieve the goal?

Binary-Search Attack
StingRay Attack

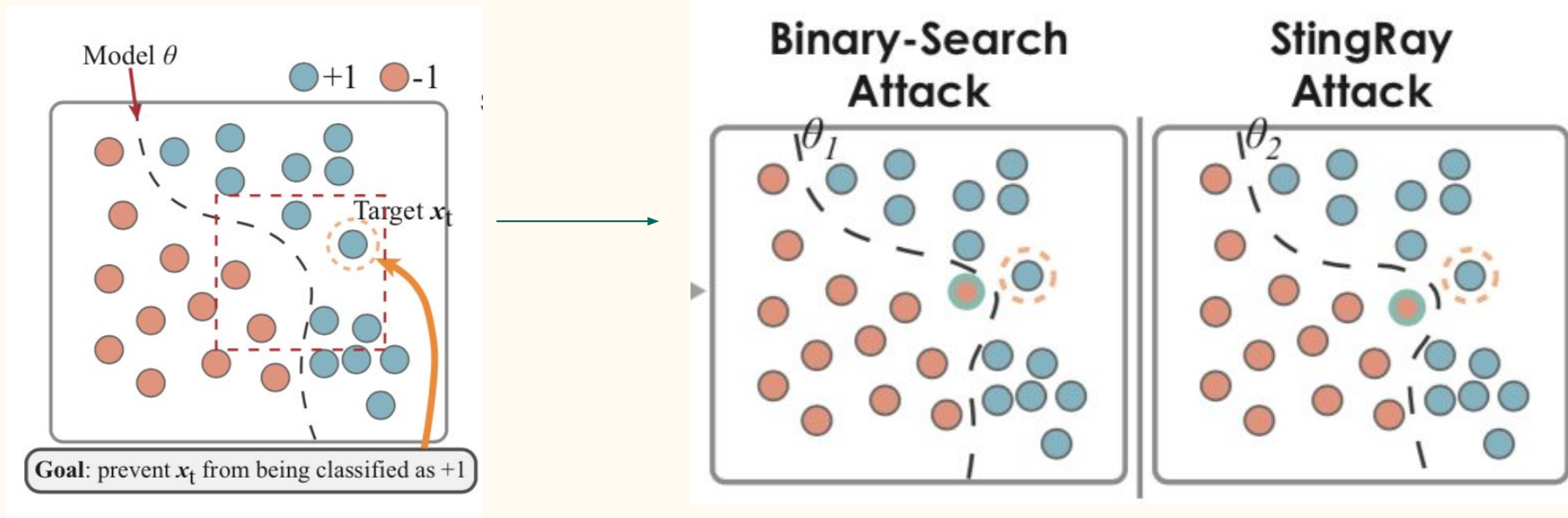
VISUAL ANALYTICS FRAMEWORK

The framework supports three main activities:

1. Vulnerability analysis
2. Attack space analysis
3. Attack results analysis

Vulnerability Analysis

- **Core idea: To change the label of the target instance**
 - **Attack algorithms: Binary-Search Attack & StingRay Attack**



Vulnerability Analysis

- **Vulnerability Measures (to explore the potential weaknesses in the model):**
 - **Decision Boundary Distances (DBD)**
 - **Minimum Cost for a Successful Attack (MCSA)**

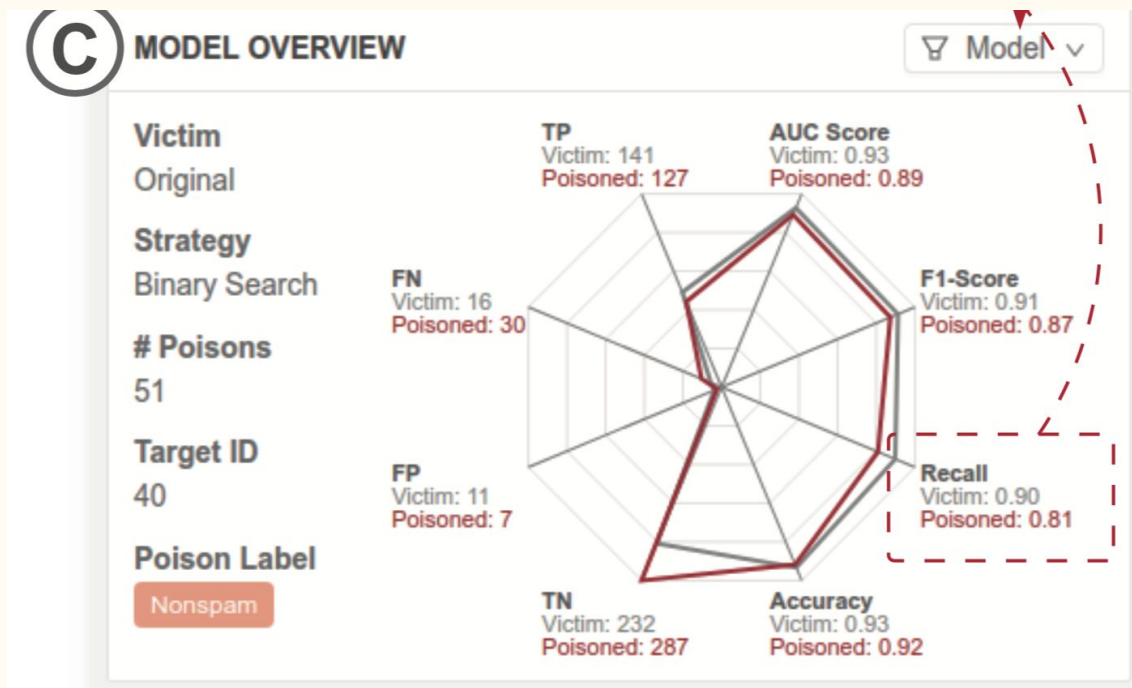
Visualizing the Attack Space

- **Vulnerability Measures:**
 - **Decision Boundary Distances (DBD)**
 - **Minimum Cost for a Successful Attack (MCSA)**
 - **Performance metrics of the poisoned model (Accuracy, Recall, etc.)**

- **Each instance in the training dataset is measured based on these vulnerability measures.**
- **[Video](#)**

Attack Detail Analysis

Model Overview



- TN
- FN
- TP
- FP

- Acc
- Recall
- F1
- ROC

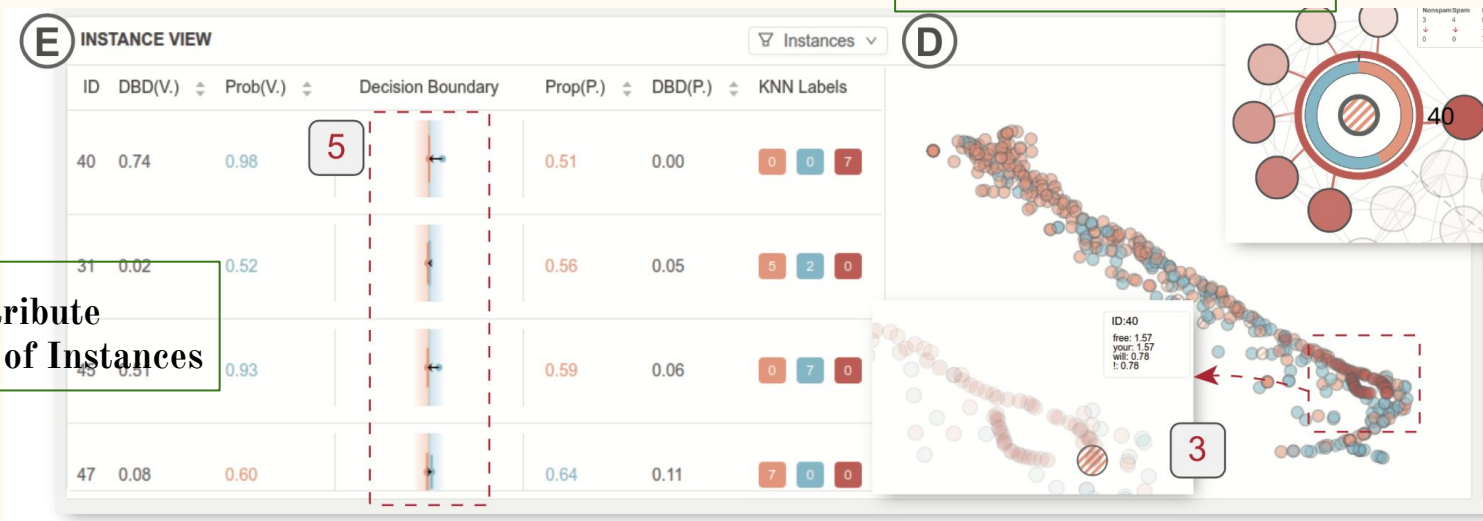
Radar Chart for Model Performances

Attack Detail Analysis

T-SNE Projection

- Overview of the data distribution

Instance View



- Instances Attribute
 - Details of Instances

- Key Attributes for Instances:
 - Decision Boundary Distances
 - Classification Probabilities
 - Labels of K-NNs

Attack Detail Analysis

Instance View

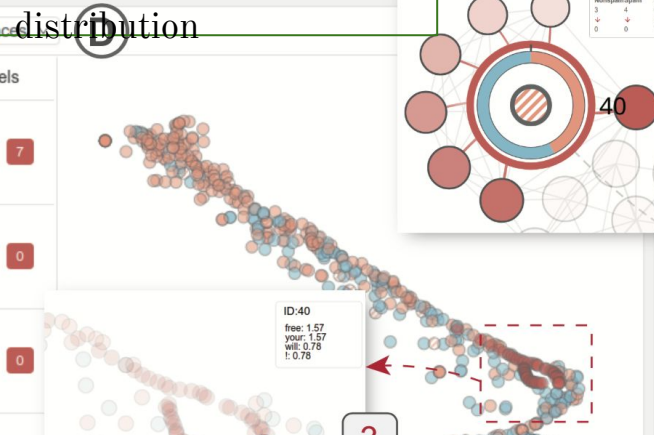
- Instances Attribute
 - Details of Instances

E INSTANCE VIEW

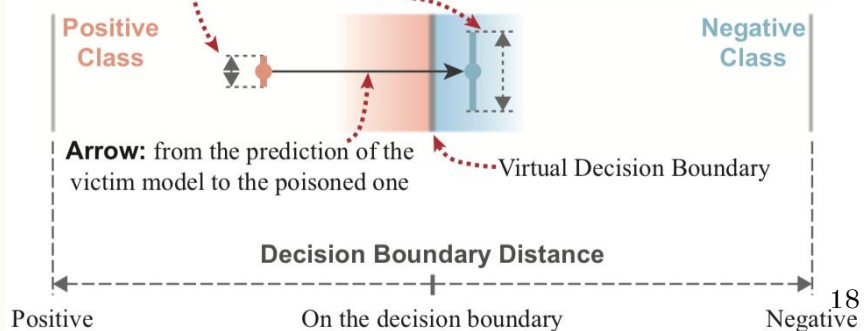
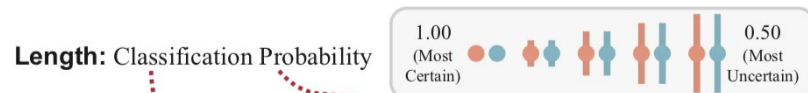
ID	DBD(V.)	Prob(V.)	Decision Boundary	Prop(P.)	DBD(P.)	KNN Labels
40	0.74	0.98		0.51	0.00	0 0 7
31	0.02	0.52		0.56	0.05	5 2 0
45	0.51	0.93		0.59	0.06	0 7 0
47	0.08	0.60		0.64	0.11	

T-SNE Projection

Overview of the data distribution



- Key Attributes for Instances:
 - Decision Boundary Distances
 - Classification Probabilities
 - Labels of K-NNs



Attack Detail Analysis

Feature View



- **Data Distributions on Features**
 - Instances in the spam/non spam classes
 - Poisoning Instances

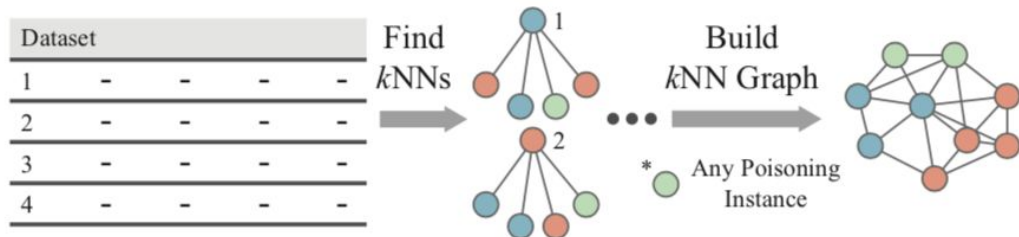
- **Feature Importance Rankings**
 - In the victim model
 - In the poisoned model
 - Differences

Attack Detail Analysis

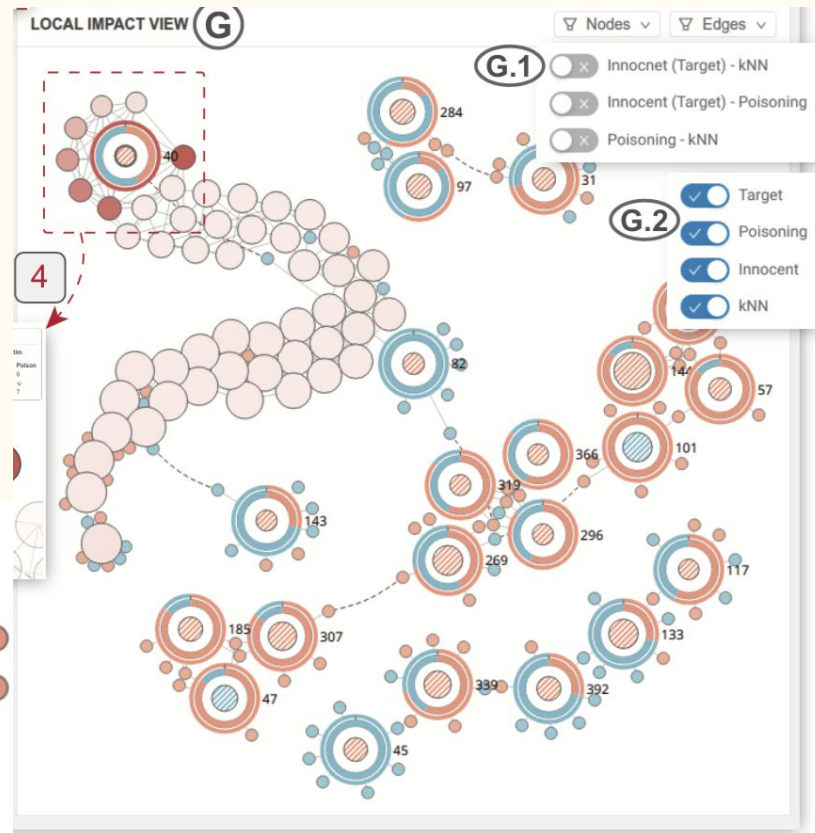
Local Impact View

- The target would be influenced either by its neighbors or the poisoned instances

KNN Graph Building



(a) kNN Graph Building

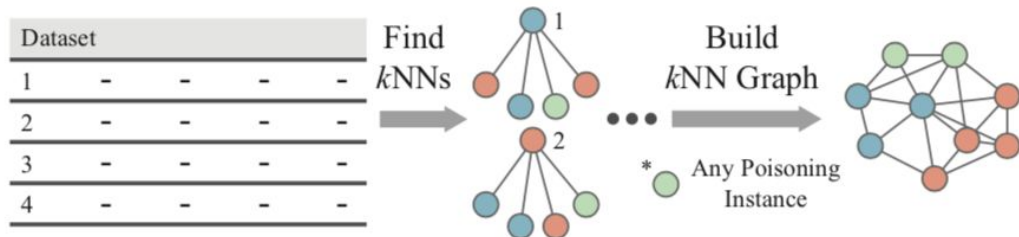


Attack Detail Analysis

Local Impact View

- The target would be influenced either by its neighbors or the poisoned instances

KNN Graph Building

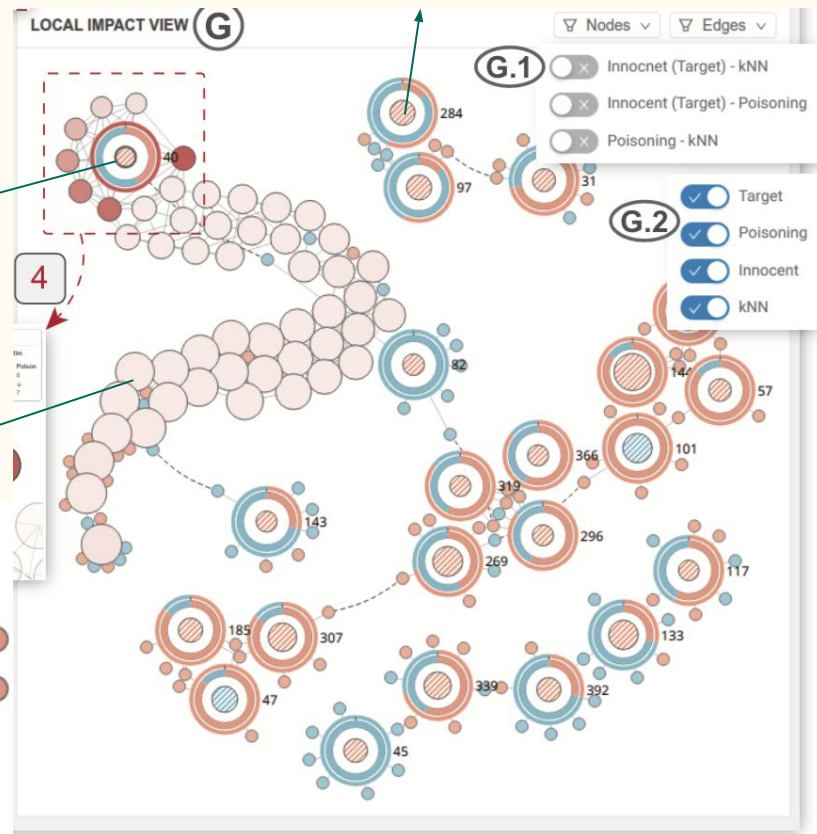


(a) kNN Graph Building

Target Instance

Poisoning Instance

Innocent Instances



Case Study

Critiques

- **Strengths**
 - **Two stage design in the interface**
 - **Very user friendly**
 - **Multi Faceted Analysis**
- **Weaknesses**
 - **Only allows to attack one instance**
 - **Speed on bigger datasets and more complicated models**
 - **Scalability (Visual design, Attack Algorithm).**
 - **Case studies are too simplified.**



SUMMARY OF ORIGINAL MODEL **A**

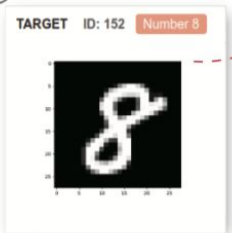
TARGET ID: 40 Spam **B**

ID	DBD	MCSA (B. Search)	Acc. (B. Search)	Recall (B. Search)	MCSA (StingRay)	Acc. (StingRay)	Recall (StingRay)	Label	Predicted
392	0.05	N/A	N/A	N/A	1	0.93 (0.00)	0.90 (0.00)	Notspam	Spam

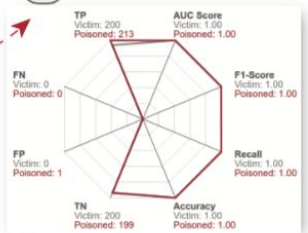
1 Data Table View

ID	DBD	MCSA (B. Search)	Acc. (B. Search)	Recall (B. Search)	MCSA (StingRay)	Acc. (StingRay)	Recall (StingRay)	Label	Predicted
152	7.41	13	1.00 (0.00)	1.00 (0.00)	14	1.00 (0.00)	1.00 (0.00)	Number 8	Number 8
126	6.19	12	1.00 (0.00)	0.99 (-0.01)	13	1.00 (0.00)	0.99 (-0.01)	Number 8	Number 8
136	5.91	12	1.00 (0.00)	1.00 (0.00)	13	1.00 (0.00)	1.00 (0.00)	Number 8	Number 8
2417	8.28	12	1.00 (0.00)	0.99 (-0.01)	13	1.00 (0.00)	0.99 (-0.01)	Number 8	Number 8
943	8.83	12	1.00 (0.00)	1.00 (0.00)	13	1.00 (0.00)	1.00 (0.00)	Number 8	Number 8

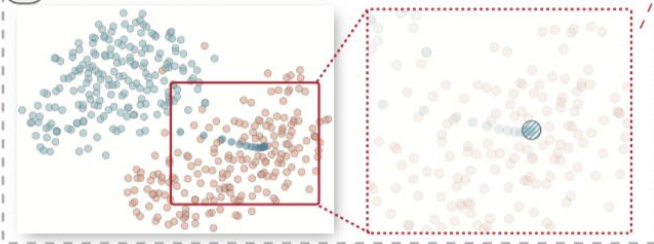
2 The Selected Target



3 Model Overview



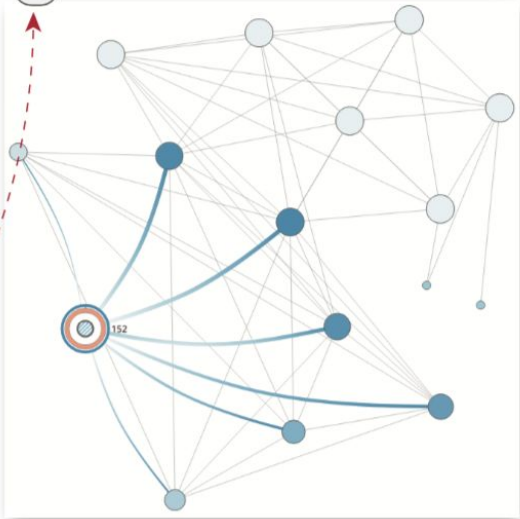
4 Projection View



6 Instance Attribute View



5 Local Impact View

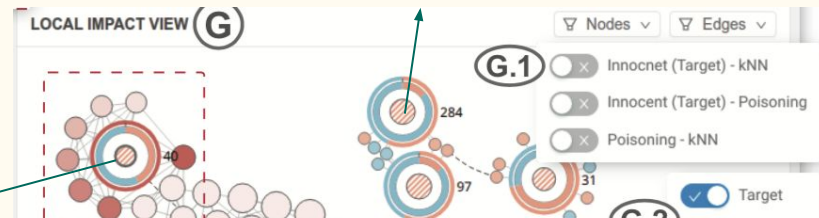


Thank you!
QA

Attack Detail Analysis

Local Impact View

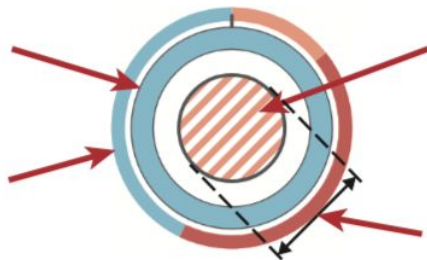
Innocent Instances



Target and Innocent Instances

Inner Ring: The class distribution of k NNs in the **victim model**

Outer Ring: The class distribution of k NNs in the **poisoned model**



Color: The predicted label

Texture: Whether the predicted label is flipped from the victim model

Size: Classification Probability

