## Derivation of Constraints from Machine Learning Models and Applications to Security and Privacy

Moreno Falaschi <u>Catuscia Palamidessi</u> Marco Romanelli







# Summary

Machine learning and explainability

#### • FOL theories from ML models

- Decision trees
- Support vector machines
- Nearest neighbors
- Neural networks

### Applications to security and privacy

- Side channel attack
- Malware
- Model inversion

#### Nowadays ML is pervasive, and performs better than humans

AlphaGo beats Go human champ



Computer out-plays humans in "doom"



Deep Net outperforms humans in image classification

IM 🔓 GENET

#### Autonomous search-and-rescue drones outperform humans



#### IBM's Watson destroys humans in jeopardy



DeepStack beats professional poker players



Deep Net beats human at recognizing traffic signs



## Generation of new examples via a GAN



Synthetic images of human faces generated using the CELEBA-HQ dataset

#### Machine learning and art



An artwork created using DeepDream, developed by researchers at Google in 2015

## On the other hand...

Original image

Perturbation A

Adversarial example



Panda 57.7% confidence





Gibbon 98.3% confidence



Temple 97% confidence





Ostrich 98% confidence

## Wrong decision can be costly and dangerous

"Autonomous car crashes, because it wrongly recognizes ..."



"Al medical diagnosis system misclassifies patient's disease ...."





The misclassification may be due to :

- insertion of special crafted noise
- wrong correlation learned in training phase





A. Esteva, et al., "Dermatologist-level classification of skin cancer with deep neural networks," Nature, 2017.

## A related problem: fairness ML can amplify bias



## Non recidivating black people twice as likely to be labelled high risk than non recidivating white people

https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

## Necessity for explanability

#### 2018





The data subject should have the right not to be subject to a decision, which may include a measure, evaluating personal aspects relating to him or her which is based solely on automated processing and which produces legal effects concerning him or her or similarly significantly affects him or her, such as automatic refusal of an online credit application or e-recruiting practices without any human intervention.

...

In any case, such processing should be subject to suitable safeguards, which should include specific information to the data subject and the right to obtain human intervention, to express his or her point of view, to obtain an explanation of the decision reached after such assessment and to challenge the decision.

# ML models as FOL constraints

- Help understanding the model
- Constraints can be manipulated
  - queries
  - check implications (derivation of new constraints)
  - satisfiability
  - ....
- Addition of new constraints
  - new knowlege
  - potential attackers
  - fairness desiderata
    - •••

# Summary

Machine learning and explainability

#### • FOL theories from ML models

- Decision trees
- Support vector machines
- Nearest neighbors
- Neural networks

## Applications to security and privacy

- Side channel attack
- Malware
- Model inversion

# **Decision trees**





(a) A binary decision tree.

$$P_{1}(X_{1}) \wedge P_{2}(X_{2}) \wedge \ldots \wedge P_{l}(X_{l})$$

$$\implies L = \ell_{1}$$

$$P_{1}(X_{1}) \wedge P_{2}(X_{2}) \wedge \ldots \wedge \neg P_{l}(X_{l})$$

$$\implies L = \ell_{2}$$

$$\vdots$$

$$\neg P_{1}(X_{1}) \wedge \neg P_{3}(X_{3}) \wedge \ldots \wedge P_{r}(X_{r})$$

$$\implies L = \ell_{t-1}$$

$$\neg P_{1}(X_{1}) \wedge \neg P_{3}(X_{3}) \wedge \ldots \wedge \neg P_{r}(X_{r})$$

$$\implies L = \ell_{t}$$

(b) The formulae derived from the decision tree.

## Support vector machines



**Constraints**  $\vec{w} \cdot \vec{X} < a \Rightarrow L = \ell_1$   $\vec{w} \cdot \vec{X} > a \Rightarrow L = \ell_2$ 

## **Nearest neighbors**



**Constraints**  $\forall \vec{w} \neq \vec{v}_1, \dots \vec{v}_k. \ (d(\vec{X}, \vec{v}_1) \leq d(\vec{X}, \vec{w}) \land \dots \land d(\vec{X}, \vec{v}_k) \leq d(\vec{X}, \vec{w})) \Rightarrow L = \ell$ 

## Neural networks

Yang et al. Explanaibility of NN through architecture constraints. IEEE TNNLS 2020

Idea: decompose the network in different sub-networks

The weights of the sub-networks are orthogonal, so that the contribution of each sub-network can be recognized in the result



Loss function

$$\begin{split} \mathcal{L}(y|m{x})) &= \mu + \sum_{j=1}^k eta_j ilde{h}_j (m{w}_j^T m{x}) + arepsilon, \ m{Constraints} \ &\sum_{i=1}^p \left| W_{ij} 
ight| \leq T_1, \ &\sum_{j=1}^k \left| eta_j 
ight| \leq T_2, \end{split}$$

 $\boldsymbol{W}^T \boldsymbol{W} = \boldsymbol{I}_k,$ 

# Summary

Machine learning and explainability

#### • FOL theories from ML models

- Decision trees
- Support vector machines
- Nearest neighbors
- Neural networks

### Applications to security and privacy

- Side channel attack
- Malware
- Model inversion

## Side channel attack

#### Example: Password checker

The program checks the enterest string bit by bit against the stored password, and fails as soon as it finds a mismatch.

Of course, if the attacker is capable to measure the time until failure, then there is a timing attack. The goal is to detect it authomatically



$$\begin{split} X_1 &= 0 \\ \implies Time = t_1 \land Result = \texttt{fail} \\ X_1 &= 1 \land X_2 = 0 \\ \implies Time = t_2 \land Result = \texttt{fail} \\ \vdots \\ X_1 &= 1 \land X_2 = 1 \land \ldots \land X_n = 0 \\ \implies Time = t_n \land Result = \texttt{fail} \\ X_1 &= 1 \land X_2 = 1 \land \ldots \land X_n = 1 \\ \implies Time = t_n \land Result = \texttt{succ} \end{split}$$

(a) The decision tree for the password example.

(b) The set of constraints derived from the decision tree for the password example.

## Side channel attack

Idea: represent the capabilities and the knowledge of the attacker by contraints

$$(X_1 = 1 \land X_2 = 1 \land \ldots \land X_m = 1) \land$$
$$(Time = t_1 \lor Time = t_2 \lor \ldots \lor Time = t_n) \land (Result = \texttt{fail} \lor Result = \texttt{succ})$$

There are n solutions for this constraint (in the context of the theory expressed in previous slide), which means that the complexity of the attack is linear.

If the time information were not present, then the number of solution (and hence the attack) would be exponential

## Malware

#### Example: a smart healtcare System



A possible attacker able to tamper with the Blood-pressure indicator

 $R-20 \leq Blood\_pressure \leq R+20 \land Diabetes\_alarm = \texttt{off}$ 

Security breach: the following scenario becomes satisfiable for any 130< r < 150

 $\begin{array}{l} Sex\doteq\texttt{female}\ ,\ Weight\doteq80\ ,\ Height\doteq165\ ,\ Age\doteq40\ ,\\ R\doteq r\ ,\ Blood\_pressure\doteq r-20\ ,\ Diabetes\_alarm\doteq\texttt{off} \end{array}$ 

Thank you !

Questions?



And, of course,

Happy Birthday, Mau!