



Cyber security systems provide, by default, a multi-agent context. Thus, one needs to consider both the aspects of cyber security and multi-agent environments to design system behaviours that provide formal bounds on security of the application. Rather than being a straightforward application of AI techniques, the cybersecurity domain also provides fresh research challenges to the AI community, as we saw in [3]. In case of MTD systems, incorporating evolution of defender configurations, attacker attacks and rewards values in the Game Theoretic framework raises the question of what can we say about optimal strategies in Repeated Games with evolutionary game metrics.

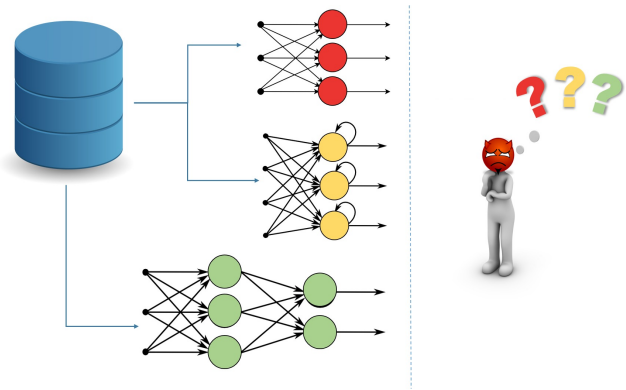
### 3. SECURITY FOR AI

In my ongoing work, I am looking at using MTD for ensuring safety of AI agents. With the use of Machine Learning algorithms in applications that affect our day to day life, we are vulnerable to attacks that seek to guide the intelligence of these approaches for malicious purposes. Consider an automated handwritten check reader at the ATM machine near you. If a malicious depositor were to put a few dots over the digit 1 so that the machine interprets it as a 9, (s)he might be able to withdraw \$900 instead of the \$100 you had planned to give him(/her). There are existing works that show that such manipulation of state-of-the-art machine learning algorithms is feasible if one can guess the type of network architecture used for such classification [5]. Although it is possible to design solutions for preventing such security compromises by reverse engineering specific attacks or incorporating adversarial examples into the train data, it is worth investigating if foundations for a general security measure is possible here.

An interesting approach would be to design multiple learners from the same testing data to keep an adversary guessing about the correct classification boundary which would make designing model-based attacks tougher. Although this Moving Target Defense approach seems related to the notion of Ensemble models, the goal of the system is to prevent adversarial samples from being misclassified as opposed to increasing classification robustness. For an MTD system to succeed in thwarting attacks, the different configurations need to have *differential immunity*. This means that adversarial samples generated for one model are ineffective (i.e. correctly classified) by all other models in the system. Existing literature has investigated such measures in the context of linear classifiers with binary labels [6], but lacks formal guarantees when these frameworks are investigated in an attacker-defender multi-agent context.

A simple idea would be to divide the data set into parts and use them to train different models for creating the configurations for the MTD framework. For the case of learning networks, as shown in [7], this idea does not provide differential immunity. On the other hand, using different network architectures to obtain the various models for creating the configurations of the MTD framework is something we are investigating at present (Figure 2). In such cases, maximizing security without sacrificing classification accuracy becomes a challenging requirement.

In the context of model-based scenarios like Markov Decision Processes or Automated Planning, it seems to be possible for an adversary to design reward shaping mechanisms for agents (without complete information) or use an agent's reward function to make them behave in an unforeseen man-



**Figure 2: Learning networks with different architectures from the same data can make it difficult for an attacker to craft malicious examples for intended mis-classification by a particular network**

ner, which may lead to dire consequences. Consider a case where a ball-catching robot is learning its reward function. An attacker learns that it has extremely high reward for catching a ball thrown at it. As the agent is deployed, the adversary can throw a ball off a cliff and tempt the robot to jump off the cliff. Notice that such instances are different from unintended consequences resulting out of bad reward function designs [8] where the robot exploits the knowledge about reward functions for itself. I plan to investigate these directions, identifying concrete problems that can lead to adversarial compromise of AI agents in these scenarios.

**Acknowledgments.** This research is supported in part by ONR grants N00014161-2892, N00014-13-1-0176, N00014-13-1-0519, N00014-15-1-2027, & the NASA grant NNX17AD06G.

### REFERENCES

- [1] M Taguinod, A Doupé, Z Zhao, and G Ahn. Toward a moving target defense for web applications. In *IEEE International Conference on Information Reuse and Integration (IRI)*, 2015, pages 510–517.
- [2] S. Vadlamudi, S. Sengupta, M. Taguinod, Z. Zhao, A. Doupé, G. Ahn, and S. Kambhampati. Moving target defense for web applications using bayesian stackelberg games. In *AAMAS*, 2016.
- [3] S. Sengupta, S. Vadlamudi, S. Kambhampati, A. Doupé, M. Taguinod, Z. Zhao, and G. Ahn. A game theoretic approach in strategy generation for moving target defense in web applications. In *AAMAS*, 2017.
- [4] A. Sinha, T.H. Nguyen, D. Kar, M. Brown, M. Tambe, and A. X. Jiang. From physical security to cyber security. *Journal of Cybersecurity*, 2016.
- [5] N Papernot, P McDaniel, S Jha, M Fredrikson, Z Celik, and A Swami. The limitations of deep learning in adversarial settings. In *Security and Privacy (EuroS&P)*, 2016 IEEE.
- [6] B Biggio, G Fumera, and F Roli. Adversarial pattern classification using multiple classifiers and randomisation. In *Joint IAPR International Workshops on Structural, Syntactic, and Statistical Pattern Recognition*, 2008.
- [7] C Szegedy, W Zaremba, I Sutskever, J Bruna, D Erhan, I Goodfellow, and R Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [8] D Amodei, C Olah, J Steinhardt, P Christiano, J Schulman, and D Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.