

3IF1135 - S curit  de l'IA

Responsables : **Pierre-Fran ois GIMENEZ**

Langues d'enseignement : **FRANCAIS**

Campus o  le cours est propos  : **CAMPUS DE RENNES**

Nombre d'heures d' tudes  l ves (HEE) : **35**

Nombre d'heures pr sentielles d'enseignement (HPE) : **15**

Ann e acad mique : **2024-2025**

Niveau avanc  : **non**

Pr sentation, objectifs g n raux du cours :

Ce module est un cours de d couverte de l'utilisation de l'IA pour la Cybers curit  ainsi que les attaques que l'on peut mener contre les m thodes bas es sur de l'IA. Le cours ne cherche pas    tre exhaustif sur l'ensemble des m thodes d'IA: au contraire, il se focalise sur les attaques contre ces syst mes, notamment celles qui ciblent les r seaux de neurones et d gradent les performances de classification, ce qui pourrait impacter les syst mes de d tection d'intrusion ou d'analyse de malware. De plus, le cours ambitionne de d crire les m thodes de l' tat de l'art qui impl mentent des m canismes de d fense dans les r seaux de neurones. Cette partie du cours sera abord e au travers de ces aspects recherche par des chercheurs du domaine. Une partie pratique permettra aux  tudiants d'impl menter des attaques et d' valuer l'efficacit  des d fenses.

Pr requis :

- Programmation en Python (cours SIP)

Plan d taill  du cours (contenu) :

- Les r seaux de neurones et de leur apprentissage (CM, 3h)
- Classification d'image par r seau de neurones convolutifs (TP, 3h)
- Attaques contre les r seaux de neurones (CM, 3h)
- D fenses contre ces attaques (CM, 3h)
- Attaque et d fense d'un classifieur (TP, 3h)

D roulement, organisation du cours :

Cours magistraux (9h)

TP (6h)

Organisation de l' valuation :

 valuation du second TP

Moyens :

Les TP utiliseront la bibliothèque pytorch (open-source).

Acquis d'apprentissage visés dans le cours :

- Développer un système de classification basé sur un réseau de neurones
- Evaluer l'impact d'attaques contre un classifieur
- Défendre un réseau de neurones contre ces attaques

Description des compétences acquises à l'issue du cours

:

C2: develop a classification system based on a neural network

C3: apply state-of-the-art attack and defense methods

Bibliographie :

- Deep Learning, Ian Goodfellow and Yoshua Bengio and Aaron Courville, MIT Press, 2016
- Adversarial Machine Learning: Attack Surfaces, Defence Mechanisms, Learning Theories in Artificial Intelligence, Aneesh Sreevallabh Chivukula, Xinghao Yang, Bo Liu, Wei Liu, Wanlei Zhou, 2023