Subject Description Form

Subject Code	DSAI5210 / COMP5577					
Subject Title	AI Security					
Credit Value	3					
Level	5					
Pre-requisite/ Co-requisite/ Exclusion	Nil					
Objectives	To equip students with a fundamental understanding of AI security and practical skills in evaluating and mitigating vulnerabilities in AI systems and leveraging AI tools and techniques to address common cybersecurity challenges. The objectives of this subject are to equip students to:					
	articulate the core concepts and principles underlying AI security;					
	identify and analyse various security threats to widely-used AI models;					
	3. comprehend and apply standard defense mechanisms to bolster the security of popular AI models;					
	4. gain proficiency in analysing and designing AI-based systems to tackle prevalent cybersecurity challenges;					
	 cultivate practical skills for implementing attack and defense strategies to build secure AI systems that address typical cybersecurity issues. 					
Intended Learning	Upon completion of the subject, students will be able to:					
Outcomes (Note 1)	Professional/academic knowledge and skills					
	a) analyse and evaluate security threats associated with different phases of AI models in deployment;					
	b) assess and execute various attacks against multiple AI models;					
	c) apply and implement common defense mechanisms to enhance the security of diverse AI models;					
	d) utilize AI techniques to address security-related challenges in common application domains.					
	Attributes for all-roundedness					
	e) acquire critical thinking and analytical skills, and improve technical writing as well as presentation skills.					

Subject Synopsis/	Topics	Number of Lectures					
Indicative Syllabus	1. Overview of AI Security	2					
(Note 2)	Introduce the fundamentals of popular AI models,	_					
	principles of security analysis, security of AI, and AI						
	for security.						
	2. Attacks on AI Model Training	2					
	Explore various attacks targeting the training phase of						
	AI models, including data poisoning attacks and						
	backdoor attacks, within domains such as computer						
	vision, network analysis, and natural language						
	processing.						
	3. Vulnerabilities of AI Models in Deployment	3					
	Analyse common threats encountered during the						
	deployment of AI models, such as evasion attacks,						
	membership inference attacks, model stealing						
	attacks, and jailbreaking attacks, across fields like						
	computer vision, network analysis, and natural						
	language processing.						
	4. Common Defense Strategies	3					
	Examine widely-used defense techniques, including						
	adversarial training, attack detection, data						
	purification, and certified robustness, to strengthen						
	the security of AI models at various stages of their						
	lifecycle.						
	5. AI for Security	3					
	Design and analyse AI-based systems to address	J					
	prevalent security challenges in various application						
	domains, such as malware detection, network						
	intrusion detection, and robust recommendation						
	systems.						
	Total	13					
	Total	13					
	TTI 11.1 1.4 1.4 1.1 1.01						
Teaching/Learning	The course will be conducted through a blend of le						
Methodology	student presentations, workshops, and a class project, e						
3.	the principles and practices of AI security. The foundational principles						
(<i>Note 3</i>)	will primarily be addressed in lectures and tutorials, while practical						
	aspects will be explored through workshops and projects. Given that AI						
	security is an evolving field characterized by rapid advancements in						
	research, students will be required to investigate selected topics and						
	present their findings during lectures. The class project will provide an						
	opportunity for students to reinforce their understanding, integrating both						
	theoretical concepts and practical skills.						

Assessment Methods in Alignment with Intended Learning Outcomes	Specific assessment methods/tasks	% weighting	Intended subject learning outcomes to be assessed (Please tick as appropriate)						
(Note 4)			a	b	c	d	e		
	1. Assignments	15	V	√	V	V			
	2. In-class Test(s)	25	V	$\sqrt{}$	V	V			
	2. Class project	20	√	$\sqrt{}$	√	V	√		
	3. Examination	40	V	√	√	V	√		
	Total	100 %			1	I			
	Explanation of the appropriateness of the assessment methods in assessing the intended learning outcomes: Assignments : assessment of the studies with respect to the understanding								
	of the relevant subject matters, including the principles, methodologies, and techniques by proving answers to the assignment questions.								
	Class project: assessment of the ability to solve real problems by using learned techniques and developing practical solutions.								
	In-class Test(s) : assessment of the level of comprehension of core concepts and the ability to apply these concepts to analyse potential threats and solutions in AI systems.								
	Examination : assessment of the overall understanding of topics.								
Student Study Effort Expected	Class contact:								
	■ Lecture					30 Hrs.			
	■ Tutorial/Student Presentation/Workshop					9 Hrs.			
	Other student study effort:								
	■ Assignments + Class Project					35 Hrs.			
	■ Self-study + Examination/Test Preparation					48 Hrs.			
	Total student study effort					122 Hrs.			
Reading List and References	1. Vorobeychik, Y., & Kantarcioglu, M. (2018). <i>Adversarial machine learning</i> . Morgan & Claypool Publishers								
	2. Miller, D. J., Xiang, Z., & Kesidis, G. (2023). <i>Adversarial Learning and Secure AI</i> . Cambridge University Press.								
	3. Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). <i>Deep learning</i> (Vol. 1, No. 2). Cambridge: MIT press.								
	4. Proceedings of IEEE Symposium on Security and Privacy								

- 5. Proceedings of USENIX Security Symposium
- 6. Proceedings of ISOC Network and Distributed System Security Symposium
- 7. Proceedings of ACM Conference on Computer and Communications Security
- 8. Proceedings of Machine Learning Research
- 9. Advances in Neural Information Processing Systems