CU_CURR501	THE CHINESE UNIVERSITY OF HONG KONG	
Page 1 of 4	Print Course Catalog Details	

Academic Org: Div of Computer Science & Engg – Subject: Al: Systems & Tech

July 25, 2024 9:18:49 AM

Course: AIST5020	Course ID: 014552	Eff Date: 2024-07-01	Crse Status: Active	Apprv. Status: Approved	[New Course]
Trustworthy Artificial Intelliger	nce 可信人工智能				

This course introduces the principles and techniques of Trustworthy Artificial Intelligence (Trustworthy AI), which aims to mitigate the potential adverse effects of AI on people and society. The course focuses on four main aspects of trustworthy AI: privacy & security, robustness, explainability, and fairness. It covers the state-of-the-art research progress in these areas, including federated learning and adversarial attacks. Algorithms, models, and systems will be covered. Moreover, the course discusses the ethical and social implications of trustworthy AI, to foster social awareness among students who would use or develop AI techniques in the future. This course is suitable for students who have some background in machine learning, probability, and linear algebra.

這門課程介紹了可信人工智能(可信AI)的原則和技術,旨在減輕AI對人們和社會可能產生的不良影響。課程專注於可信AI的四個主要方面:隱私與安全、韌性、可解釋性和 公平性。課程涵蓋了這些領域的最新研究進展,包括聯邦學習和對抗性攻擊。課程將涵蓋算法、模型和系統。此外,課程討論了可信AI的道德和社會影響,以培養學生對未來 使用或開發AI技術的社會意識。該課程適合具有一定機器學習、概率和線性代數背景的學生。

Grade Descriptor: A

EXCELLENT – exceptionally good performance and far exceeding expectation in all or most of the course learning outcomes; demonstration of superior understanding of the subject matter, the ability to analyze problems and apply extensive knowledge, and skillful use of concepts and materials to derive proper solutions.

有關等級說明的資料,請參閱英文版本。

В

GOOD – good performance in all course learning outcomes and exceeding expectation in some of them; demonstration of good understanding of the subject matter and the ability to use proper concepts and materials to solve most of the problems encountered.

有關等級說明的資料,請參閱英文版本。

С

FAIR – adequate performance and meeting expectation in all course learning outcomes; demonstration of adequate understanding of the subject matter and the ability to solve simple problems.

有關等級說明的資料,請參閱英文版本。

D

MARGINAL – performance barely meets the expectation in the essential course learning outcomes; demonstration of partial understanding of the subject matter and the ability to solve simple problems.

有關等級說明的資料,請參閱英文版本。

F

FAILURE – performance does not meet the expectation in the essential course learning outcomes; demonstration of serious deficiencies and the need to retake the course.

有關等級說明的資料,請參閱英文版本。

Equivalent Offering:	
Units:	3 (Min) / 3 (Max) / 3 (Acad Progress)
Grading Basis:	Graded
Repeat for Credit:	Ν
Multiple Enroll:	Ν
Course Attributes:	MSc Computer Science
	MPhil-PhD Computer Sci & Erg

Topics:

COURSE OUTCOMES

Learning Outcomes:

At the end of the course of studies, students will be able to:

1. Understand and explain the main concepts and techniques of trustworthy AI and its applications in various domains;

2. Compare and contrast the advantages and disadvantages of different trustworthy AI techniques;

3. Implement and evaluate different types of trustworthy AI techniques using popular frameworks such as PyTorch or TensorFlow;

CU_CURR501	THE CHI	NESE UNIVERSITY OF HONG KONG	July 25, 2024
Page 3 of 4		Print Course Catalog Details	9:18:49 AM
	4. Discuss the ethical and social in	nplications of trustworthy AI, in terms of privacy, security, robustn	ess, explainability, fairness,
	etc.; 5. Organize and conduct a group r	roject for state of the art research	
	5. Organize and conduct a group p		
Course Syllabus:			
	Week 1. Introduction to trustworthy	/ AI	
	Week 2. Privacy & Security: federa	ited learning basics	
	Week 3. Privacy & Security: recons	struction attacks and inference attacks	
	Week 4. Privacy & Security: anony	mization, noise injection, SMPC, and homomorphic encryption	
	Week 5. Robustness: outliers and	adversarial attacks	
	Week 6. Robustness: adversarial o	letection and defense	
	Week 8 Explainability: post-boc ex	rplainability	
	Week 9. Fairness: individual fairne	ss and group fairness	
	Week 10. Fairness: fairness-aware	e federated learning	
	Week 11. Ethical and social implication	ations of trustworthy Al	
	Week 12. Group presentation I		
	Week 13. Group presentation II		
Assessment Type:	Homework or assignment	• 30%	
	Presentation	: 30%	
	Project	: 40%	
Feedback for Evaluation:			
	1. Results of homework and assigr	nments.	
	2. Course evaluation and question	naire.	
	3. Reflection of teachers.		
	4. Question-and-answer sessions	during class.	
	5. Student consultation during offic	e hours or online.	
Required Readings:			
	1 Trustworthy Machine Learning	Kush P. Varshney, Independently published, 2022	
	2 A Survey of Trustworthy Federa	ted Learning with Perspectives on Security Robustness and Priv	vacy. Yifei Zhang, Dun Zeng
		tes _earling many eropeentes on occarty, robusticess, and r m	

CU_CURR501	
Page 4 of 4	

	Jinglong Luo, Zenglin Xu, Irwin King. ArXiv. 2023.
Recommended Re	eadings: 1. Practical Adversarial Robustness in Deep Learning: Problems and Solutions. Pin-Yu Chen, Sayak Paul. CVPR 2021.
	OFFERINGS
1. AIST5020	Acad Organization=CSEGV; Acad Career=RPG
	COMPONENTS
	LEC : Size=30: Final Exam=N: Contact=3
	TUT : Size=30; Final Exam=N; Contact=1
	ENROLMENT REQUIREMENTS
1. AIS15020	Enrollment Requirement Group: For students in MSc Computer Science; or
	For students in MPhil-PhD Computer Science & Engineering
	New Enrollment Requirement(s):
	Other Requirement = For students in MSc Computer Science; or
	For students in MPhil-PhD Computer Science & Engineering
	Additional Information
	eLearning hrs for blended cls 0 VTL-Onsite face-to-face hrs 0 VTL-Online synch. hrs 0 VTL-Online asynch. hrs 0 No. of micro-modules 0

< E N D O F R E P O R T >