# YUZHE MA

Microsoft Corporation, Microsoft Azure AI

yuzhema@microsoft.com & ma234@wisc.edu · 6082133287 · Google Scholar

## Education

**University of Wisconsin-Madison** *09/2021*
Ph.D. in Computer Sciences, Minor in Statistics
Advisor: Professor Xiaojin (Jerry) Zhu

**University of Wisconsin-Madison** *05/2018*
M.S. in Computer Sciences

**Huazhong University of Science and Technology** *06/2016*
B.E. in Computer Science and Technology
Advisor: Professor Kun He

## Research Interests

My research interest lies broadly in machine learning and artificial intelligence. My current research focuses primarily on Natural Language Processing (NLP) and its applications. In particular, I have worked on building modern NLP techniques such as distillation on Azure AI platform, benchmarking the performance and stability of Large Language Models (LLMs) across various tasks and datasets, and delivering products based on LLMs to Microsoft internal teams. Another topic of my research is adversarial sequential decision making (check a tutorial here). Within this topic, I focus on analyzing adversarial vulnerability/robustness of models in typical sequential decision making scenarios, including multi-armed bandit, reinforcement learning, optimal control systems, and multi-agent game-theoretical learning scenarios. In addition to the above two topics, my previous research has also touched unsupervised learning (e.g., dimensionality reduction, clustering, social networks etc.), deep neural networks, machine teaching, and differential privacy in machine learning.

## Work Experience

**Senior Data & Applied Scientist, Microsoft** *Since 10/2021*

○ Developed a library called Babel that enabled internal product teams to perform end-to-end experiments with state-of-the-art Large Language Models (LLMs) such as GPT-3, GPT-3.5, ChatGPT, and GPT-4. The platform supports both model finetuning and few-shot inference on customer data.

○ Built a distillation pipeline for compressing large LLMs into more efficient and frugal smaller ones. The distilled model are demonstrated to achieve higher accuracy than directly finetuning the student model.

○ Developed a benchmarking platform that supports fast evaluation of LLMs on more than 100 datasets. The library is able to benchmark both the quality and the throughput performance (e.g., RPS) of the models. The library also applies to multimodal models such as GPT-Vision.

**Research Intern, IBM Research** *06-08/2020*

○ Built a two-stage machine learning model based on Cycle Generative Adversarial Network (Cycle-GAN) and object detection techniques to identify buildings on the satellite imagery data of American cities.

○ Developed an iterative training procedure based on the object detection algorithm YOLO to augment the labels in the training data. The model increased the total amount of labeled buildings from 20K to 80K.

○ Applied Cycle-GAN to transform imagery data into rasterized maps.

**Applied Scientist Intern, Amazon** *06-09/2019*

- Developed a student identification model using the Gradient Boosted Tree (GBT) algorithm.
- Carried out an end-to-end machine learning pipeline, including data acquisition, model training, hyperparameter tuning, post-processing of predictions, and model testing.
- Evaluated the model performance on Amazon Prime data and achieved 85% accuracy.

### Research Intern, Symantec Research Labs (NortonLifeLock)          *05-08/2018*

- Proposed a federated machine teaching framework that coordinates the training process of multiple computer nodes to jointly teach a desired target model, while preserving the privacy of local datasets on each node during the communication and information sharing process.
- The work was published in the International Joint Conference on Neural Networks (**IJCNN**).

### Research Scholar, Cornell University          *07-08/2015*

- Studied manifold learning and developed techniques for nonlinear dimensionality reduction.
- Investigated theoretical properties of dimensionality reduction methods when the data has an underlying manifold structure. The work was published in Theoretical Computer Science (**TCS**).
- Advisor: Professor Kun He & Professor John E. Hopcroft.

### Research Assistant, John Hopcroft Lab @ Huazhong University          *02/2013-05/2016*

- Worked on unsupervised machine learning, including clustering and dimensionality reduction.
- Developed a method to detect overlapping communities in social networks. The work was published in the National Conference on Theoretical Computer Science (**NCTCS**).
- Developed a dimensionality reduction algorithm that can preserve both global and local data structure. The work was published in The International Frontiers of Algorithmics Workshop (**FAW**).

## Refereed Publications

($\alpha$-$\beta$) indicates that the authors are listed in alphabetical order.

[1]. Jing Cui, Yufei Han, **Yuzhe Ma**, Jianbin Jiao, and Junge Zhang. BadRL: Sparse Targeted Backdoor Attack Against Reinforcement Learning. In The 38th AAAI Conference on Artificial Intelligence (**AAAI**), 2024.

[2]. **Yuzhe Ma**, and Zhijin Zhou. Adversarial Attacks on Adversarial Bandits. In The 11th International Conference on Learning Representations (**ICLR**), 2023. (Spotlight, acceptance rate 5.65%)

[3]. **Yuzhe Ma**, Young Wu, and Xiaojin Zhu. Game Redesign in No-regret Game Playing. In The 31th International Joint Conference on Artificial Intelligence (**IJCAI**), 2022.

[4]. **Yuzhe Ma**, Young Wu, and Xiaojin Zhu. Game Redesign in No-regret Game Playing. In The NeurIPS Learning in Presence of Strategic Behavior Workshop (**NeurIPS-LPSB**), 2021.

[5]. Yun-Shiuan Chuang, Xuezhou Zhang, **Yuzhe Ma**, Mark K. Ho, Joseph L. Austerweil, and Xiaojin Zhu. Using Machine Teaching to Investigate Human Assumptions when Teaching Reinforcement Learners. In The 43rd Annual Meeting of the Cognitive Science Society (**CogSci**), 2021.

[6]. **Yuzhe Ma**, Jon Sharp, Ruizhe Wang, Earlence Fernandes, and Xiaojin Zhu. Demo: Sequential Attacks on Kalman Filter-based Forward Collision Warning Systems. In The NDSS Automotive and Autonomous Vehicle Security Workshop (**NDSS-AutoSec**), 2021. (Demo)

[7]. **Yuzhe Ma**, Jon Sharp, Ruizhe Wang, Earlence Fernandes, and Xiaojin Zhu. Sequential Attacks on Kalman Filter-based Forward Collision Warning Systems. In The 35th AAAI Conference on Artificial Intelligence (**AAAI**), 2021.

[8]. Xuezhou Zhang, Shubham Bharti, **Yuzhe Ma**, Adish Singla, and Xiaojin Zhu. The Sample Complexity of Teaching by Reinforcement on Q-learning. In The 35th AAAI Conference on Artificial Intelligence (**AAAI**), 2021.

[9]. Xuezhou Zhang, **Yuzhe Ma**, Adish Singla. Task-agnostic Exploration in Reinforcement Learning. In The 34th Conference on Neural Information Processing Systems (**NeurIPS**), 2020.

[10]. Xuezhou Zhang, **Yuzhe Ma**, Adish Singla, and Xiaojin Zhu. Adaptive Reward-Poisoning Attacks against Reinforcement Learning. In The 37th International Conference on Machine Learning (**ICML**), 2020.

[11]. **Yuzhe Ma**, Xuezhou Zhang, Wen Sun, Xiaojin Zhu. Policy Poisoning in Batch Reinforcement Learning and Control. In The 33rd Conference on Neural Information Processing Systems (**NeurIPS**), 2019.

[12]. **Yuzhe Ma**, Xiaojin Zhu, and Justin Hsu. Data Poisoning against Differentially-Private Learners: Attacks and Defenses. In The 28th International Joint Conference on Artificial Intelligence (**IJCAI**), 2019.

[13]. Yufei Han, **Yuzhe Ma**, Chris Gates, Kevin Roundy, and Yun Shen. Collaborative and Privacy-Preserving Machine Teaching via Consensus Optimization. In The International Joint Conference on Neural Networks (**IJCNN**), 2019.

[14]. Kwang-Sung Jun, Lihong Li, **Yuzhe Ma**, and Xiaojin Zhu. Adversarial Attacks on Stochastic Bandits. In The 32nd Conference on Neural Information Processing Systems (**NeurIPS**), 2018. (**α-β**)

[15]. **Yuzhe Ma**, Kwang-Sung Jun, Lihong Li, and Xiaojin Zhu. Data Poisoning Attacks in Contextual Bandits. In The 9th Conference on Decision and Game Theory for Security (**GameSec**), 2018.

[16]. Ayon Sen, Scott Alfeld, Xuezhou Zhang, Ara Vartanian, **Yuzhe Ma**, and Xiaojin Zhu. Training Set Camouflage. In The 9th Conference on Decision and Game Theory for Security (**GameSec**), 2018.

[17]. **Yuzhe Ma**, Robert Nowak, Philippe Rigollet, Xuezhou Zhang, and Xiaojin Zhu. Teacher Improves Learning by Selecting a Training Subset. In The 21st International Conference on Artificial Intelligence and Statistics (**AISTATS**), 2018.

[18]. **Yuzhe Ma**, Kun He, John Hopcroft, and Pan Shi. Neighbourhood-Preserving Dimension Reduction via Localised Multidimensional Scaling. In Theoretical Computer Science (**TCS**), 2017.

[19]. **Yuzhe Ma**, Kun He, John Hopcroft, and Pan Shi. Nonlinear Dimension Reduction by Local Multidimensional Scaling. In The 10th International Frontiers of Algorithmics Workshop (**FAW**), 2016.

[20]. **Yuzhe Ma**, Kun He, Leihua Qin, and Yan Wang. A Primary Research on Overlapping Community Detection. In The 32nd National Conference on Theoretical Computer Science (**NCTCS**), 2014.

## Other Publications

[1]. Amin Saied, **Yuzhe Ma**, Mercer Chen, Gilsinia Lopez, and Ali Mahmoudzadeh. BabelBench: Large Language Model Benchmarking in AI Platform. In The Microsoft Machine Learning, AI & Data Science Conference (**MLADS**), 2024. (**Internal**)

[2]. Clarisse Simoes, Lars Liden, Swadheen Shukla, Subho Mukherjee, Yu Wang, Luciano Del Corro, **Yuzhe Ma**, and Amin Saied. Building an Automated Distillation Pipeline to Democratize Large-Scale AI. In The Microsoft Machine Learning, AI & Data Science Conference (**MLADS**), 2023. (**Internal**)

[3]. **Yuzhe Ma**. Adversarial Attacks in Sequential Decision Making and Control. Published by The University of Wisconsin-Madison, 2021. (**PhD Thesis**)

[4]. Yufei Han, **Yuzhe Ma**, Chris Gates, Kevin Roundy, Yun Shen. Systems and Methods for Preventing Decentralized Malware Attacks, U.S. Patent, 11,025,666, 2021.

## Honors and Awards

| | |
|---|---|
| Student Travel Award, NeurIPS | *2019* |
| Top 50% Reviewer, NeurIPS | *2019* |
| Student Travel Award, GameSec | *2018* |

| | |
|---|---|
| Honorarium Award, GameSec Special Track | *2018* |
| Student Travel Award, AISTATS | *2018* |
| UW CS Summer Research Award | *2017* |
| Outstanding Bachelor Thesis Award | *2016* |
| CCF Outstanding Undergraduate Award (only 100 awardees nationwide) | *2015* |
| Academic Excellence Scholarship | *2014* |
| Outstanding Student Leader Award | *2014* |
| China National Scholarship Award | *2013* |
| Merit Student Award | *2013* |

## Academic Service

Program Committee: ECAI23, ACML23, ACML22, AAAI22, ACML21, AAAI21, ACML20, AAAI20, ACML19, AAAI19

Conference Reviewer: CogSci24, ICML24, ICLR24, AISTATS24, NLDL24, NeurIPS23, ICML23, CogSci23, AISTATS23, ICLR23, NeurIPS22, CogSci22, ICML22, AISTATS22, ICLR22, NeurIPS21, ICML21, AISTATS21, ICLR21, NeurIPS20, ICML20, AISTATS20, NeurIPS19, ICML19, AISTATS19

Journal Reviewer: TMLR, TON, TPAMI, IEEE Access, Machine Learning

Student Volunteer: AISTATS18

## Student Contest

| | |
|---|---|
| **Student Cluster Competition of Super-Computing Conference 2014** | *2014* |
| Ranked No. 5 for overall and No. 3 for Linpack | |
| **Fifth National Undergraduate Mathematical Contest of China** | *2013* |
| First Prize in Hubei Province | |
| **Chinese Mathematical Olympiad in Senior** | *2011* |
| First Prize in Jiangsu Province | |
| **Chinese Physics Olympiad in Senior** | *2011* |
| Second Prize in Jiangsu Province | |

## Teaching Experience

| | |
|---|---|
| **University of Wisconsin-Madison** | *09/2016-01/2017* |
| Teaching Assistant - Introduction to Computer Engineering (CS252) | |

## Talks & Presentations

| | |
|---|---|
| Adversarial Attacks on Adversarial Bandits. | *02/19/2023* |
| *Microsoft Azure AI Group* | |
| Adversarial Example Attacks and Defenses in NLP. | *09/06/2022* |
| *Microsoft Azure AI Group* | |
| Adversarial Machine Learning in Sequential Decision Making. | *11/05/2021* |
| *Microsoft Azure AI Group* | |
| Data Poisoning against Differentially-Private Learners: Attacks and Defenses. | *08/14/2019* |
| *IJCAI 2019 at China, Macau.* | |
| Machine Teaching Theory and Its Applications. | *12/26/2018* |
| *Huazhong University of Science and Technology.* | |

Data Poisoning Attacks in Contextual Bandits. *10/30/2018*
*GameSec 2018 at University of Washington.*

Teacher Improves Learning by Selecting a Training Subset. *04/20/2018*
*1st IFDS Student Workshop at University of Wisconsin-Madison.*

## Skills & Expertise

**Programming Skills**: Python, Pytorch, SQL, Matlab, C, C++, R, AMPL, Verilog

**Machine Learning**: Adversarial Machine Learning, Multi-armed Bandit, Reinforcement Learning, Game Theory, Natural Language Processing (NLP), Advanced Driver Assistance Systems (ADAS), Differential Privacy, Dimensionality Reduction, Machine Teaching