# University of Tehran
# School of Electrical and Computer Engineering

| Course | Interpretable AI | | |
|---|---|---|---|
| **Course type, level, credit** | Optional | Graduate | 3 units |
| **Field, Major** | Electrical and Computer Engineering | All majors | |
| **Co-requisite(s)** | - | | |
| **Prerequisite(s)** | | | |
| **Prerequisite by topic** | - | | |
| **Goals** | The use of machine learning in practical settings, is different than in a lab setting. In practical settings, several issues appear, that do not generally appear in lab settings. Some of these practical issues include: Robustness, Domain Generalization, Interpretability, Explainability, Safety, Security, Fairness, Transparency, privacy and regulatory. <br><br> The goal of this course is to equip students with the knowledge and experience of these practical problems of machine learning and allow them to make more informed decision in their design of machine learning systems. | | |
| **Outcome** | Upon successful completion of the course, the students will gain a knowledge and understanding about the practical issues of using AI in practice. | | |
| **Topics** | 1- Robustness and Generalization: <br>    - General AI <br>        - Bias-Variance trade-off and regularization <br>        - Inductive biases vs. generalization <br>        - Thought process to choose inductive biases <br>        - How to use inductive biases in a model <br>        - Self-supervision and its useful inductive biases <br>        - Multi-task learning and its emerging properties <br>        - The emergence of deduction (AI) from induction (learning) <br>        - Learning with loss function vs learning with rules <br><br>    - Generalization vs robustness <br>        - Generalization to different domains <br>        - Robustness to input noise in a given domain <br>        - trade-offs between robustness and generalization <br><br>    - Learning with noisy inputs <br>        - Effects of noise on generalization | | |

| | |
|---|---|
| | - data augmentation and associated inductive biases<br><br>- Learning with low-quality supervision<br>    - Robustness to noisy labels<br>    - Evaluation given noisy labels<br>    - Sampling, rare events and black swan problem<br>    - Sparse and low-entropy labels<br>    - Anomaly detection<br><br>2- Interpretability and Explainability:<br>    - Prescriptive models (models that make decisions for us)<br>        - Why predictive models need some causal understanding<br>        - How prescriptive models make decisions affecting humans<br>        - Tragedies in prescriptive machine learning<br>        - Why interpretability matters in prescriptive models<br><br>    - Interpretable models<br>        - Linear models, decision trees, etc. and how to interpret them.<br>        - The relation between model interpretability vs causal analysis<br>        - Trade-off between interpretability and prediction accuracy<br>        - Practical cases where interpretability becomes more important<br>than accuracy<br><br>    - Explaining the workings of machine learning models<br>        - Explainability vs interpretability<br>        - Importance of explaination in advancing machine learning<br>        - Model-dependent vs model-agnostic methods<br>        - Example-based methods<br>        - Global and local explaination techniques<br>        - Understanding the workings of deep models<br><br>3- Safety and Security:<br>    - Security attacks in machine learning:<br>        - poisoning of a dataset<br>        - backdoor attacks and trojan attacks<br>        - open-access models vs closed-access models<br>        - Security attacks in generative models and language models<br>        - Watermarking<br><br>    - Defense against attacks:<br>        - Trojan detection and characterization<br>        - Input/output monitoring<br><br>    - Safety and reliability:<br>        - Supervised/unsupervised Anomaly detection<br>        - Out of distribution detection<br>        - Accumulation of error in autoregressive models<br>        - Safety in prescriptive models<br>        - How to quantify reliability (domain shift, augmentation,<br>adversarial)<br><br>4- Fairness, Transparency, ethics and Privacy<br>    - Fairness and justice<br>        - How to define fairness in mathematical terms<br>        - Controversies in expressing defining justice mathematically<br>        - How to measure fairness/bias<br>        - Use of interpretable models to analyze fairness<br>        - Ways to enforce fairness in machine learning models |

| | |
|---|---|
| | - Transparency and Opqueness in decision making<br>      - In what contexts decision making process should be transparent<br>      - Importance of interpretability in transparency<br>      - Transparency regulations<br><br>- Ethics of AI<br>      - How AI affects the job market, democracy, environment<br>      - Who is legally liable if an AI causes an injury<br>      - The use of chat bots and their impact.<br>      - Fiar use of data<br>      - Regulatory bodies of AI |
| **Required software** | - Python, PyTorch |
| **Assignments** | This course will cover four projects. (Details below) |
| **Projects** | Students are expected to work on four course projects throughout the semester. Projects must be completed individually. Projects cover:<br>- Robustness<br>- Interpretability<br>- Machine Learning Security<br>- Fairness |
| **Grading** | Projects:           50 %<br>Midterm exam:     20 %<br>Final exam:        30 % |
| **Textbook(s)** | https://christophm.github.io/interpretable-ml-book/<br><br>https://fairmlbook.org/<br><br>https://fairmlclass.github.io/<br><br>http://www.math.ku.dk/~peters/elements.html |
| **Further readings** | https://www.fatml.org/<br><br>https://facctconference.org/ |
| **Prepared by** | Mohammad Amin Sadeghi, Mostafa Tavassolipour |

Prepared by: Mohammad Amin Sadghi, Last revision 20.10.1401