## CS 281 Ethics of AI – Reading list 2025

This is a list of suggested (but not mandatory) readings for the course sorted by topic. The list is not meant to be static, and we will be updating entries to this list as the course progresses.

## Fairness & Bias

- (textbook) FAIRNESS AND MACHINE LEARNING Limitations and Opportunities
- <u>Fairness</u>, <u>Equality</u>, <u>and Power in Algorithmic Decision-Making</u>
- Equality of opportunity in supervised learning
- Fairness Through Awareness
- o Delayed Impact of Fair Machine Learning
- Learning Fair Representations
- Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings
- <u>Learning controllable Fair Representations</u>
- FACT: A Diagnostic for Group Fairness Trade-offs
- Right Decisions from Wrong Predictions: A Mechanism Design Alternative to Individual Calibration
- Retiring Adult: New Datasets for Fair Machine Learning
- The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning
- o On Fairness and Calibration
- Calibration for the (Computationally-Identifiable) Masses
- Predicting Good Probabilities With Supervised Learning

## Explainability

- (textbook) Interpretable Machine Learning A Guide for Making Black Box Models Explainable
- "Why Should I Trust You?": Explaining the Predictions of Any Classifier
- A Unified Approach to Interpreting Model Predictions
- Actionable Recourse in Linear Classification
- Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV)
- Understanding Black-box Predictions via Influence Functions
- Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving Models

- o Sanity Checks for Saliency Maps
- <u>Learning to Explain: An Information-Theoretic Perspective on Model</u> Interpretation
- o Interpretable Machine Learning: Moving From Mythos to Diagnostics
- Stop Explaining Black Box Machine Learning Models for High Stakes
  Decisions and Use Interpretable Models Instead

## Privacy

- o (textbook) The Algorithmic Foundations of Differential Privacy
- o Privacy as Contextual Integrity
- o Deep Learning with Differential Privacy
- o Machine Unlearning
- o A Taxonomy of Privacy