

Machine Learning Theory (CS 6783)

Tu-Th 1:00 to 2:15 PM

Uris Library, 2B02

Online: Zoom link on course website

<https://www.cs.cornell.edu/courses/cs6783/2021fa>

Instructor : Karthik Sridharan

ABOUT THE COURSE

- No exams !
- 4 assignments that count towards your grades (40%)
- One term project (60%)

ASSIGNMENTS

- You are allowed a total of 7 days of late submission (across the 4 assignments)
- Assignments submissions via cms in pdf format

TERM PROJECT

- One research style project
- Literature survey due mid semester (will announce date later)
- Project report due at the end
- A short presentation for fun towards end of semester
- Projects can be done in group sizes of at most 2

PRE-REQUISITES

- This is a theory course and will be heavy on math!
- Basic probability theory
- Basics of algorithms and analysis
- Introductory level machine learning course
- *Mathematical maturity, comfortable reading/writing formal mathematical proofs.*

Lets get started ...

WHAT IS MACHINE LEARNING

Use **past** observations to **automatically learn** to make better predictions/decisions in the **future**.

WHERE IS IT USED ?

Recommendation Systems

NETFLIX

Browse Task Profile KIDS DVDs

PLAY

Titles, People, Genres

Karthik

House of Cards 2013-2014 TV-MA 2 Seasons

NETFLIX ORIGINAL
HOUSE OF CARDS

Bad, for a greater good.

Season 2 of this acclaimed original thriller series earned a total of 13 Emmy Award nominations including Outstanding Drama Series, Outstanding Lead Actor nominee Kevin Spacey stars as ruthless, cunning Congressman Francis Underwood, who will stop at nothing to conquer the halls of power in Washington D.C. His secret weapon: his gorgeous, ambitious, and equally conniving wife Claire (Outstanding Lead Actress nominee Robin Wright).

Directors' Commentary Available

Watch Season 1 of this Emmy-winning series with exclusive scene-by-scene audio commentary from directors including David Fincher and Joel Schumacher.

Genres: [TV Shows](#), [TV Dramas](#)

This show is: [Witty](#), [Cerebral](#), [Dark](#)

★★★★★

Our best guess for Karthik: 4.9 stars

Average of 4,007,827 ratings: 4.5 stars

+ My List

WHERE IS IT USED ?

Pedestrian Detection



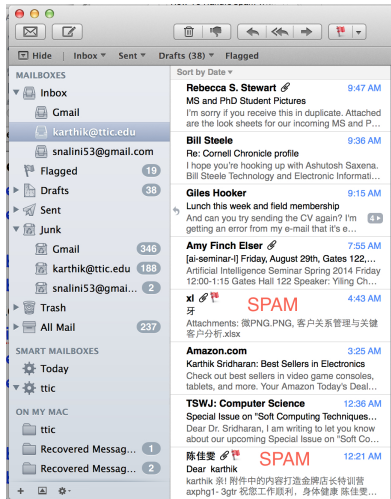
WHERE IS IT USED ?

Market Predictions



WHERE IS IT USED ?

Spam Classification



WHERE IS IT USED ?

- Online advertising (improving click through rates)
- Climate/weather prediction
- Text categorization
- Unsupervised clustering (of articles ...)
- ...

WHAT IS LEARNING THEORY

- “Develops a rigorous and quantitative understanding of machine learning problems”
 - Formalize the machine learning problem setup
 - Understand fundamental limitations: how well can we hope to learn, how many samples do we need, how do we collect them, how computationally efficient can we be
 - Develop algorithm design tools that under the said formalism provides “provable guarantees” on performance
- Understand and analyze existing approaches mathematically, when do they succeed and when do they fail, when should you use which algorithm

WHAT IS MACHINE LEARNING THEORY

- How do we formalize machine learning problems
- Right framework for right problems (Eg. online , statistical)
- How do we pick the right model to use and what are the tradeoffs between various models
- How many instances do we need to see to learn to given accuracy
- How do we design learning algorithms with provable guarantees on performance
- *Computational learning theory : which problems are efficiently learnable*

FORMALIZING LEARNING PROBLEMS

- How is data generated ?
- How do we measure performance or success ?
- Where do we place our prior assumption or model assumptions ?

FORMALIZING LEARNING PROBLEMS

- How is data generated ?
- How do we measure performance or success ?
- Where do we place our prior assumption or model assumptions ?
- *What we observe ?*

OUTLINE OF TOPICS

- Learning problem and frameworks, settings, minimax rates
- Statistical learning theory
 - Probably Approximately Correct (PAC) and Agnostic PAC frameworks
 - Empirical Risk Minimization, Uniform convergence, Empirical process theory
 - Bound on learning rates: MDL bounds, PAC Bayes theorem, Rademacher complexity, VC dimension, covering numbers, fat-shattering dimension
 - Supervised learning : necessary and sufficient conditions for learnability
- Online learning theory
 - Sequential minimax and value of online learning game
 - Regret bounds: Sequential Rademacher complexity, Littlestone dimension, sequential covering numbers, sequential fat-shattering dimension
 - Online supervised learning : necessary & sufficient conditions for learnability
- Algorithms for online convex optimization: Exponential weights algorithm, strong convexity, exp-concavity and rates, Online mirror descent
- Deriving generic learning algorithms : relaxations, random play-outs
- If time permits, uses of learning theory results in optimization, approximation algorithms, perhaps a bit of bandits, ...

THREE SCENARIOS

- 1 Recognizing which species of animal is it from images (or gambling at a roulette table)
- 2 Investing a penny a day based on stock market expert advice
- 3 Investing a pot of gold (a fixed amount) for n days based on stock market expert advice

How do we formalize each one?

SCENARIO I: STATISTICAL LEARNING

SCENARIO II: ONLINE LEARNING

SCENARIO III: ONLINE LEARNING WITH STATES (POLICY REGRET)

LEARNING PROBLEM : BASIC NOTATION

- Input space/ feature space : \mathcal{X}
(Eg. bag-of-words, n-grams, vector of grey-scale values, user-movie pair to rate)
- Output space/ label space \mathcal{Y}
(Eg. $\{\pm 1\}$, $[K]$, \mathbb{R} -valued output, structured output)
- Loss function : $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$
(Eg. 0-1 loss $\ell(y', y) = \mathbf{1}\{y' \neq y\}$, sq-loss $\ell(y', y) = (y - y')^2$), absolute loss $\ell(y', y) = |y - y'|$
Measures performance/cost per instance (inaccuracy of prediction/ cost of decision).
- Model class/Hypothesis class $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$
(Eg. $\mathcal{F} = \{x \mapsto f^T x : \|f\|_2 \leq 1\}$, $\mathcal{F} = \{x \mapsto \text{sign}(f^T x)\}$)

A FIRST EXAMPLE: ONLINE BINARY CLASSIFICATION

$\mathcal{Y} = \{\pm 1\}$, $\ell(y', y) = \mathbf{1}\{y' \neq y\}$, $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ is a finite set of models

SCENARIO ONE: IID SAMPLES

- At round t learner receives sample $x_t \sim \mathbf{D}_X$ (drawn iid from fixed distribution)
- Learner makes prediction $\hat{y}_t \in \{\pm 1\}$
- $y_t = f^*(x_t)$ is revealed where $f^* \in \mathcal{F}$ is unknown to learner

What is our strategy? How well does it do, how many mistakes?

SCENARIO ONE: IID SAMPLES

SCENARIO TWO: BINARY CLASSIFICATION ARBITRARY SEQUENCE

- At round t learner receives instance $x_t \in \mathcal{X}$ chosen arbitrarily
- Learner makes prediction $\hat{y}_t \in \{\pm 1\}$
- $y_t = f^*(x_t)$ is revealed where $f^* \in \mathcal{F}$ is unknown to learner

What is our strategy? How well does it do, how many mistakes?

SCENARIO TWO: BINARY CLASSIFICATION

ARBITRARY

SNEEK PEEK

- Different learning frameworks (more formally)
- No Free Lunch Theorems
- Minimax rates for various setting/problems
- Comparing the various settings