

CS 598: Machine Learning for Sys, Networks, and Security

Gang Wang

University of Illinois Urbana-Champaign

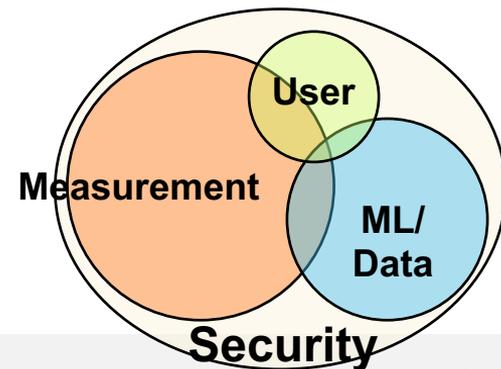
Fall 2020

Gang Wang – About Me



<https://gangw.cs.illinois.edu>

- 2006 - 2010: B.E. Tsinghua
 - 2010 - 2016: Ph.D. UC Santa Barbara
 - 2016 - 2019: Assistant Professor at Virginia Tech
 - 2019 – Present: Assistant Professor at University of Illinois
-
- **Research:** Security and Privacy, Data Mining, Measurement, Mobile Net
 - **Main publication venues:** USENIX Security, CCS, IEEE SP, NDSS, IMC, WWW, CHI, AAAI
 - 2019: IMWUT Distinguished Paper Award
 - 2018: CCS Outstanding Paper
 - 2018: NSF CAREER Award
 - 2017: Google Faculty Award
 - 2013: SIGMETRICS Best Paper



Get to Know You

- CS, ECE, or other departments?
- PhD, MS, or Undergraduate?
- ML or Sys/Net/Security background?
 - Have run/implement some ML algorithm before?
- How many years in your current program?



Class Logistics

- Class meeting:
 - WF 12:30 PM - 01:45 PM (Central Time)
 - <https://illinois.zoom.us/j/92312908719?pwd=TENGZTgzV0kyZzZITIFUM2dqUndHUT09>
 - If you cannot use zoom over browser, try to use the zoom client app
 - Meeting ID: 923 1290 8719
 - Password: 244874
- Office hours
 - I will stay longer after the class meeting for any questions
 - Contact me to make an appointment
- If any student needs special accommodations because of a disability, please contact me in the first week of class



What We Want to Achieve in This Class

This is not a conventional machine learning class = You need security/sys background

- Understanding the recent applications of machine learning (ML) in security and networked systems
- Identifying new problems / solutions in security/system domains
- Identifying the limitations and misuse of machine learning
- Applying ML to your own area in creative ways?



Paper Reading

This is not a conventional Machine Learning Class



- ML for defense
 - Spam, phishing, intrusion, malware, code analysis, deception detection
- ML for attack
 - Voice, image, deepfake, password guessing, Tor privacy attack
- NLP vs. security
 - Privacy policy, vulnerability reports
- Attacking and securing ML
 - Poisoning, evasion, trojanning/backdoor (briefly, system aspect)

Paper Reading (Cont')

This is not a conventional Machine Learning Class



- ML explanations
 - Post-hoc, human-machine collaboration, vulnerabilities and limitations
- ML for networking
 - Routing, scheduling, protocol design, user applications
- ML vs. software engineering, robot sensors
 - Debugging, testing, control (brief)

Main sources of papers: **USENIX Security**, **CCS**, **IEEE SP**, **NDSS**,
CVPR, **AAAI**, **NeurIPS**, **SIGCOMM**

Course Website

- <https://gangw.web.illinois.edu/class/cs598/>
- Easy to find through my homepage
- Your go-to place
 - Project description
 - Paper list and schedule
 - Announcement
 - Class policies

CS 598: Machine Learning for Sys, Networks, and Security

[Home](#) | [Project](#) | [Piazza](#)

Instructor	Gang Wang (gangw@illinois.edu)
Time/Location	WF 12:30 PM - 01:45 PM. Zoom information in this Google Doc (need to use your illinois Google App to view)
Office Hour	By Appointment

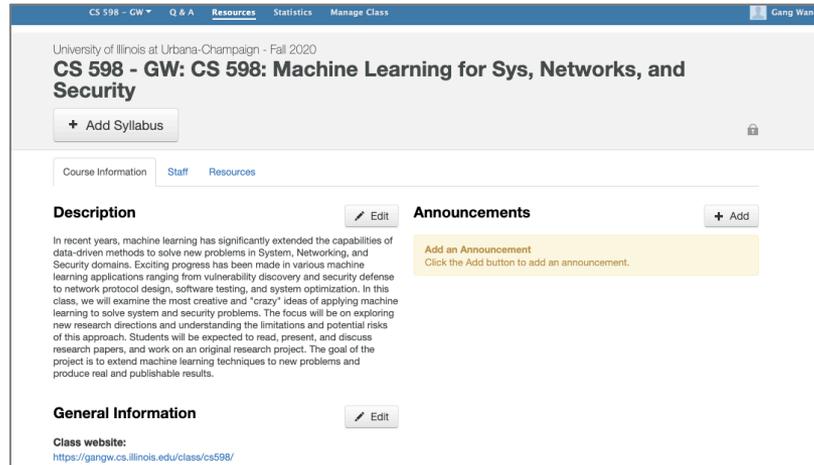
Class Description

In recent years, machine learning has significantly extended the capabilities of data-driven methods to solve new problems in System, Networking, and Security domains. Exciting progress has been made in various machine learning applications ranging from vulnerability discovery and security defense to network protocol design, software testing, and system optimization. In this class, we will examine the most creative and "crazy" ideas of applying machine learning to solve system and security problems. The focus will be on exploring new research directions and understanding the limitations and potential risks of this approach. Students will be expected to read, present, and discuss research papers, and work on an original research project. The goal of the project is to extend machine learning techniques to new problems and produce real and publishable results.



Piazza

- <https://piazza.com/class/kdl3b8mcs3p2j2>
- **Important:** receive announcement, find teammate, QA,
- **Sign up now** to claim paper to present



The screenshot shows the Piazza interface for a course. At the top, there are navigation tabs: CS 598 - GW, Q & A, Resources, Statistics, and Manage Class. The course title is "University of Illinois at Urbana-Champaign - Fall 2020 CS 598 - GW: CS 598: Machine Learning for Sys, Networks, and Security". Below the title is a "+ Add Syllabus" button. The main content area has tabs for "Course Information", "Staff", and "Resources". Under "Course Information", there are sections for "Description" and "General Information", each with an "Edit" button. The "Description" section contains a paragraph about machine learning applications in system and security domains. The "Announcements" section has an "Add" button and a yellow box with the text "Add an Announcement" and "Click the Add button to add an announcement." The "General Information" section lists the "Class website" as <https://gangw.cs.illinois.edu/class/cs598/>.



Expected Work (Reading, Participation)

- **Attend all the zoom lectures**
 - Email me ahead of time if you cannot attend (e.g., paper deadline, sick leave)
- **Paper reading**
 - 2 papers per class, be original, post reviews to Piazza (one thread per paper)
 - If you are the first, you get to create the thread and summarize the paper.
 - Other students should avoid repeating the same arguments/comments that the previous reviews have already covered
 - Post before 11:50 AM (CT) on the day of class
- **Paper presentation**
 - Build your own slides, don't directly use the authors' slides stack (20-25 mins)
 - Discuss the **key idea**, main contributions/findings, strengths and limitations
 - Discuss what you would like to do differently or the new ideas you have
 - Sign up close on Sunday 11:59 pm. **Instructor will sign the slot manually afterward**



How to Select a Paper to Present

- Sign up class Piazza (class website → piazza)
- Find the question: “select the paper you want to present”

CS 598: Machine Learning for Sys, Networks, and Security

[Home](#) | [Project](#) | [Piazza](#)

Instructor	Gang Wang (gangw@illinois.edu)
Time/Location	WF 12:30 PM - 01:45 PM. Zoom information in this Google Doc (need to use your illinois Google App to view)
Office Hour	By Appointment

Class Description

In recent years, machine learning has significantly extended the capabilities of data-driven methods to solve new problems in System, Networking, and Security domains. Exciting progress has been made in various machine learning applications ranging from vulnerability discovery and security defense to network protocol design, software testing, and system optimization. In this class, we will examine the most creative and “crazy” ideas of applying machine learning to solve system and security problems. The focus will be on exploring new research directions and understanding the limitations and potential risks of this approach. Students will be expected to read, present, and discuss research papers, and work on their original research project. The goal of the project is to extend machine



Expected Work (Project)

- **Team project**
 - <https://gangw.web.illinois.edu/class/cs598/project.html>
 - Goal: original research, publishable results
 - Fall 2019 → three USENIX Security submissions
- **Project topics: ML + X**
 - X should be roughly within “security”, “system”, and “networks”
 - We will talk about some suggested “X” next class
 - Do something you are familiar with
 - Better if it is aligned with your own research
 - Group of 3 students (same grade for all group members)

Proposal → Mid-term → Status update → Final presentation / report



Expected Work: Project

- **Proposal**
 - 1 page, describe the problem, background, and your idea
 - **Publishable idea**, talk to me before starting (after class, office hours)
- **Midterm presentation + report**
 - The research problem, your idea, and preliminary results
- **Status update slides**
 - Some slides to report the project progress
- **Final presentation**
 - Presenting your idea, your approaches, key findings and results
- **Project report**
 - 6-page paper



Project Timeline

[Home](#) | [Project](#) | [Piazza](#)

- 09/16/2019: Project proposal due (1-page)
- 10/07/2019: mid-term presentation (focusing on your idea)
- 10/09/2019: mid-term report due (abstract, background, literature survey)
- 11/06/2019: progress update slides due (dataset, experiment design, preliminary result)
- 12/09/2019: final presentation (focusing on your results and findings)
- 12/15/2019: final report due (the full paper)

*All deadlines are 11:59 PM (CT) of the specific date.



Grading

- Sum up the points: x out of 100
- Convert x to letter grade:
 - [0-60] F, [60-62] D-, [63-66] D, [67-69] D+, [70-72] C-, [73-76] C, [77-79] C+, [80-82] B-, [83-86] B, [87-89] B+, [90-92] A-, [93-100] A.
- Do not curve the grades

Class attendance and participation	5%
Paper reviews	25%
Paper presentation in class	10%
Project: proposal	10%
Project: midterm presentation	10%
Project: final presentation	15%
Project midterm report + status update	10%
Final project report	15%



Policies

- Late policy
 - All the deadlines are hard deadlines
 - Late submission score = $0.5 \times (\text{your raw score})$, if delay < 72 hours (3 days)
 - Late submission score = 0, if delay \geq 72 hours (3 days)
 - This does not apply to final project report (must be on time)
- Academic Integrity
 - Make proper citations for other people's ideas, tools, code, datasets
 - Write your own and project reports
 - **Paper presentation:** you should not copy / directly use the authors' slides



In-Classroom Policies

Active Participation

You are encouraged to turn on your video (not required if it is incontinent) and use your audio to ask and answer questions



ToDos This/Next Week

This week:

- Sign up for Piazza
- Sign up for the presentation slot
- Find a team for projects
- Submit your review + comment to Piazza before the Friday class

Next week:

- Work on project proposal
- Talk to me /email me your project idea



A Brief Introduction



What is Machine Learning?

The complexity in traditional computer programming is in the code. In machine learning, algorithms (programs) are in principle simple and the complexity (structure) is in the data. Is there a way that we can automatically learn that structure? That is what is at the heart of machine learning.

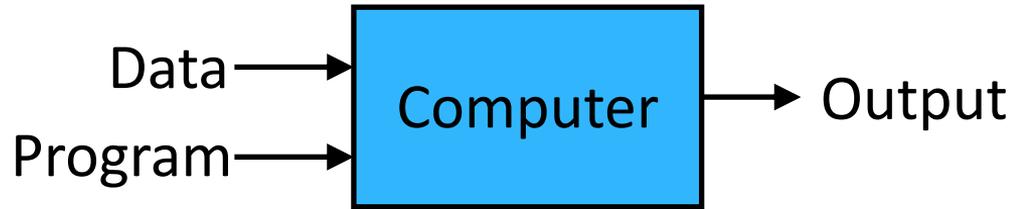
-- Andrew Ng

Machine learning is the about the construction of systems that can learn from data.

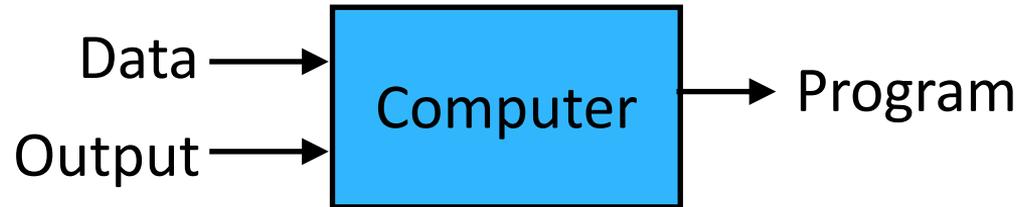


What is Machine Learning?

Traditional Programming



Machine Learning



ML is a form of Induction

- Given examples of a function $(x, f(x))$
 - *Don't explicitly know f*
 - Labeled training data set, i.e., $f(x)$
 - Training set will be noisy, i.e., $(x, (f(x) + \epsilon))$
- Predict function **$f(x)$** for new examples x
 - Discrimination/Prediction (Regression): $f(x)$ continuous
 - Classification: $f(x)$ discrete
 - Estimation: $f(x) = P(Y = c|x)$ for some class c



When To Use Machine Learning?

- When patterns exist in our data
 - Even if we don't know what they are
- We can not pin down the functional relationships mathematically
 - Else we would just code up the algorithm
- When we have lots of (unlabeled) data
 - Labeled training sets harder to come by
 - Data is of a high-dimension
 - Want to discover lower-dimension representations



Examples of ML Problems

- Pattern Recognition
 - Facial identities or facial expressions
 - Handwritten or spoken words (e.g., Siri)
 - Medical images
 - Network traffic
- Pattern Generation
 - Generating images
 - Motion sequences



Examples of ML Problems

- Outlier Detection
 - Unusual patterns in data center networks
 - Unusual sequences of credit card transactions
 - Unusual patterns of sensor data from a nuclear power plant
- Prediction
 - Future stock prices or currency exchange rates
 - Network events



What To Learn

- Learning is a procedure that estimates the parameters for a statistical model that can perform certain tasks
- Different types of learning considered here
 - Supervised
 - Unsupervised
 - Semi-supervised learning
 - Reinforcement learning



Supervised Learning

- The desired output (label) of given inputs is well known
 - You are given the “answer” (label) in the training set
 - Training data is set of $(x^{(i)}, y^{(i)})$ pairs, $x^{(i)}$ is the input example, $y^{(i)}$ is the label
- There are many supervised learning algorithms
 - Examples: Decision Trees, Ensembles (Bagging, Boosting, Random Forests, ...), k-NN, Linear Regression, Naive Bayes, Logistic Regression, Support Vector Machines (and other Large Margin Classifiers), Deep Neural Networks



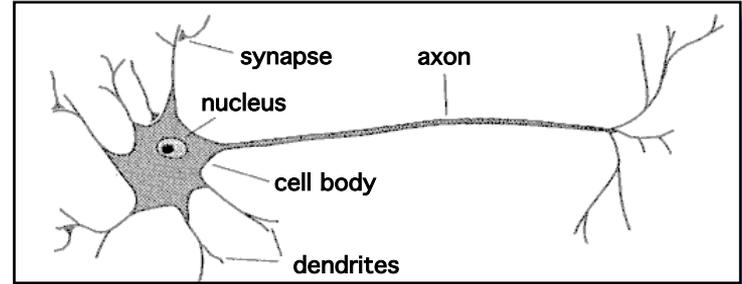
Unsupervised Learning

- Basic idea: discover unknown structure in input data
 - No need for labeled data
- Data clustering and dimension reduction
 - More generally: find the relationships/structure in the data set
- Learning algorithms include (again, many algorithms)
 - K-Means Clustering
 - Auto-encoders/deep neural networks
 - Restricted Boltzmann Machines
 - Sparse Encoders



Biological Inspiration: Neurons

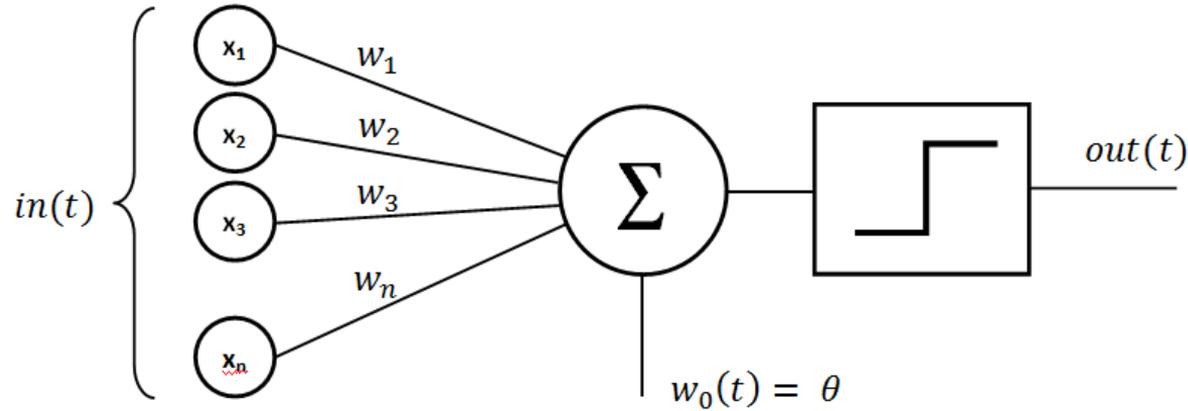
- A neuron has
 - Branching input (dendrites)
 - Branching output (the axon)



- Information moves from dendrites to axon via the cell body
- Axon connects to dendrites via synapses
 - Synapses vary in strength
 - Synapses may be excitatory or inhibitory

Basic Perceptron

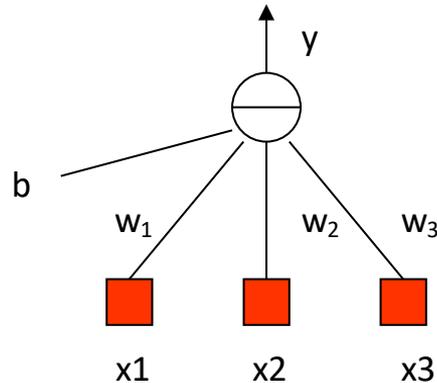
(Rosenblatt, 1950s and early 60s)



$$O = \left\{ \begin{array}{l} 1 : \left(\sum_i w_i x_i \right) + b > 0 \\ 0 : otherwise \end{array} \right\}$$

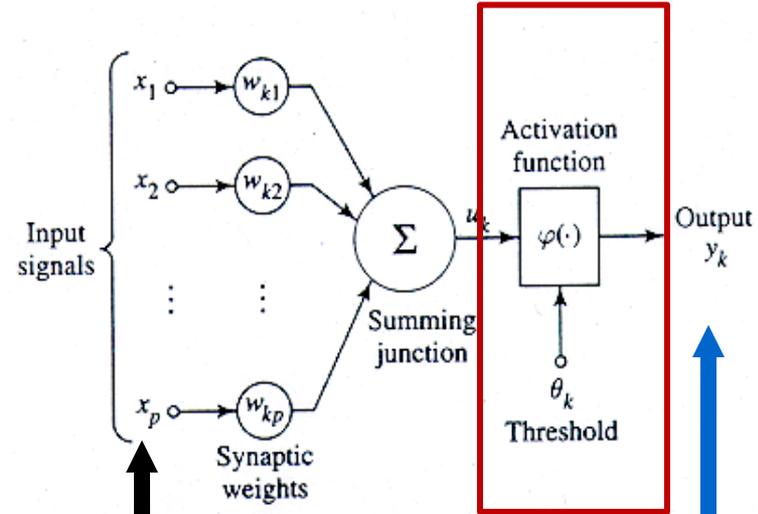
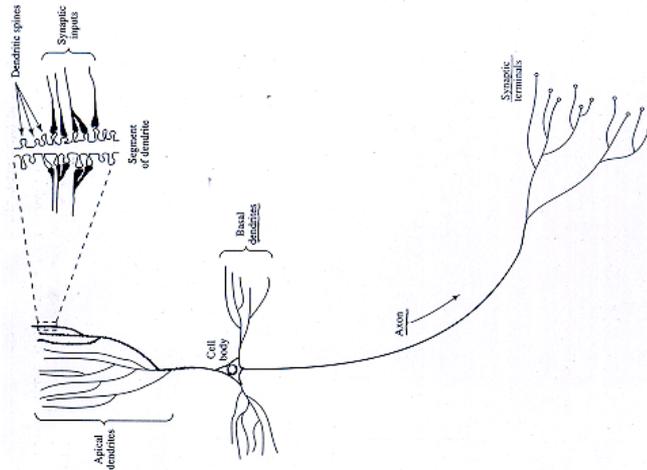
Artificial Neuron?

- An Artificial Neuron (AN) is a non-linear parameterized function with restricted output range

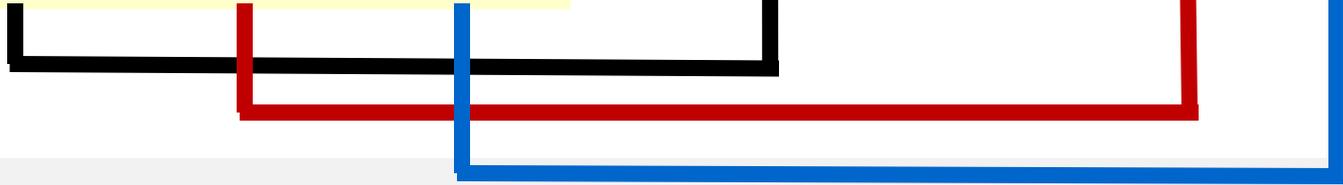


$$y = f\left(b + \sum_{i=1}^{n-1} w_i x_i\right)$$

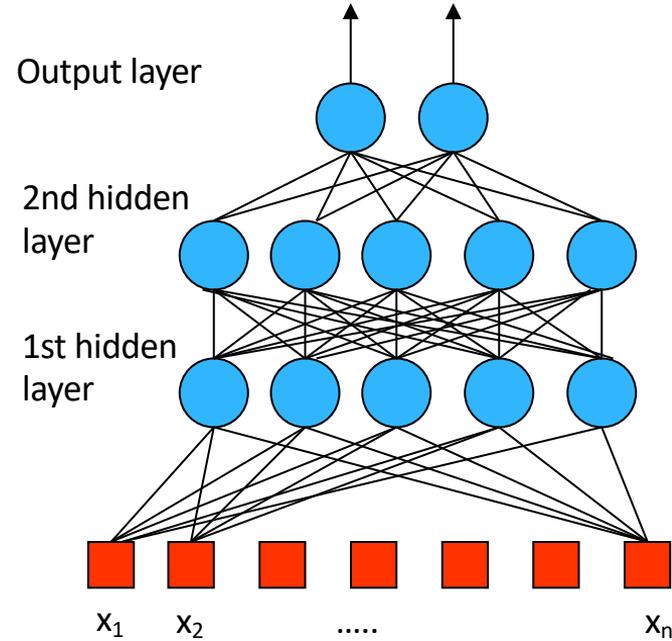
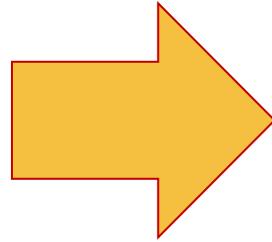
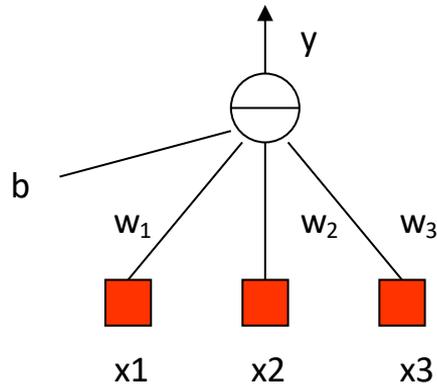
Mapping to Biological Neurons



Dendrite Cell Body Axon

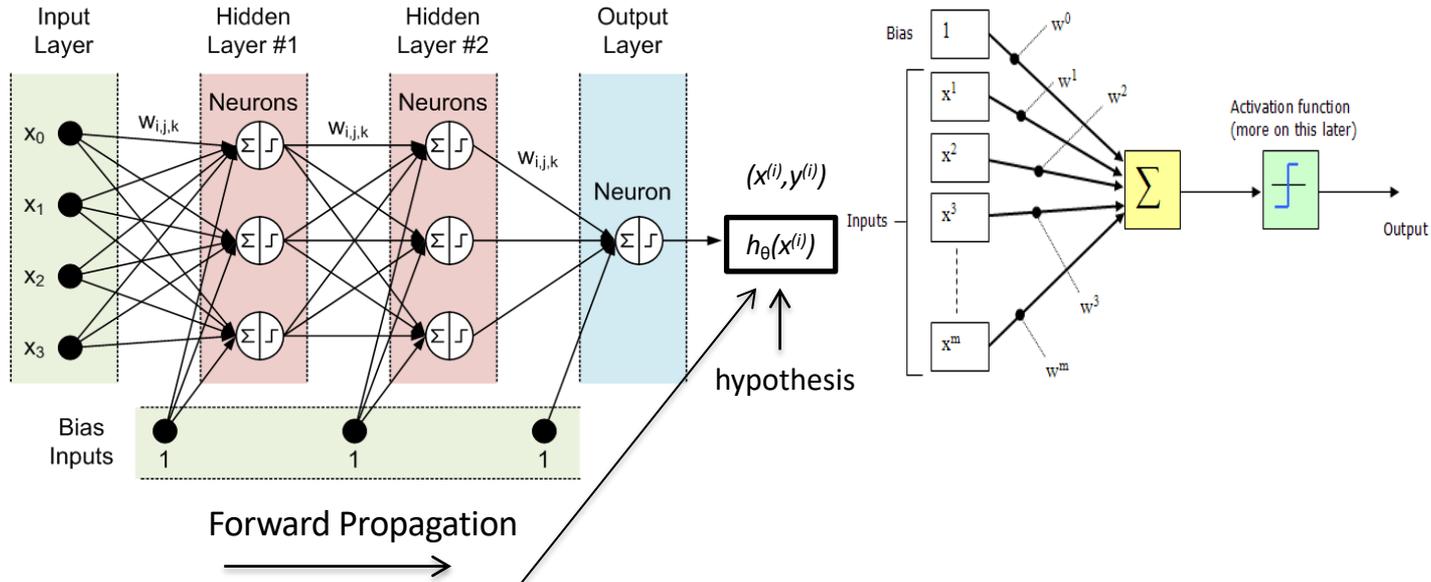


Neural Networks



input data, or features, are n dimensional

Deep Feed Forward Neural Nets



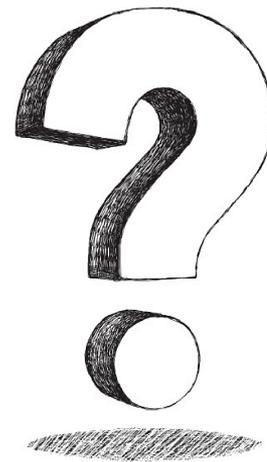
$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



Learning is the adjusting of the weights $w_{i,j}$ such that the cost function $J(\theta)$ is minimized
Simple learning procedure: *Back Propagation* (of the error signal)

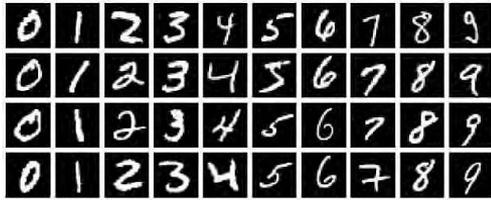
ML vs. Security/Sys/Nets

- New ways of solving classic problems
- Solutions to new problems
- Problems of ML based solutions
- Creating new problems (e.g., new attacks)



Example: Pattern Recognition

- Helpful apps: digitizing books, pedestrian detection, navigation



- New attacks: solving CAPTCHA, more powerful malicious bots



Pattern Generation/Synthesis

- Generating new faces that never exist in the world
- Style transfer (digital arts)



<https://hardikbansal.github.io/CycleGANBlog/>



- New attacks: phishing, public deception (fake news)

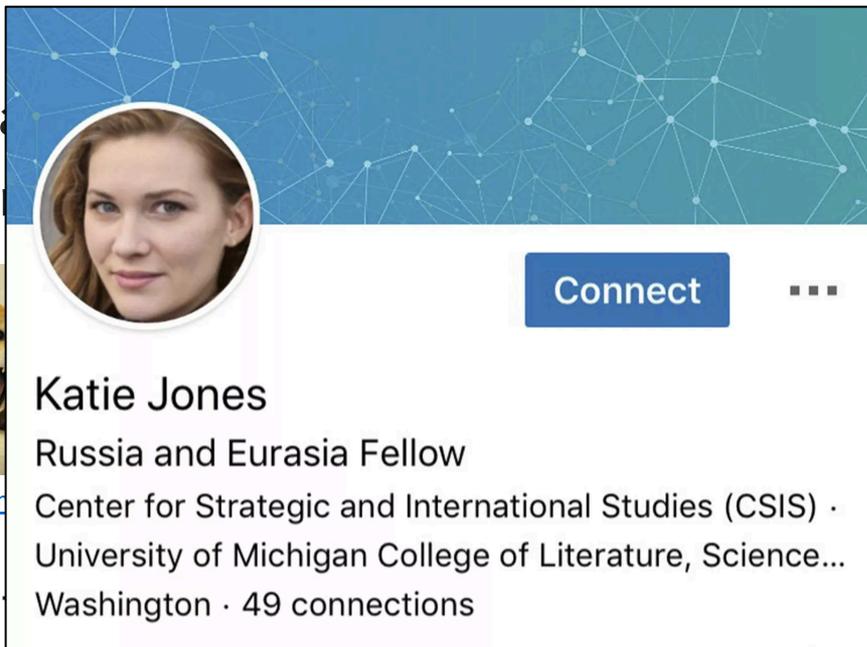
Pattern Generation

- Genera
- Style t



<https://hardikbar>

- New a



Katie Jones
Russia and Eurasia Fellow
Center for Strategic and International Studies (CSIS) ·
University of Michigan College of Literature, Science...
Washington · 49 connections

the world



<http://thispersondoesnotexist.com/>



(fake news)

Experts: Spy used AI-generated face to connect with targets

By RAPHAEL SATTER June 13, 2019

Deep Fake

- Generating/modifying image/videos by mimicking a “reference”



Reference

Our Result

<https://gizmodo.com/deepfake-videos-are-getting-impossibly-good-1826759848>



Bill Hader



Arnold
Schwarzenegger

<https://www.youtube.com/watch?v=bPhUhypV27w>

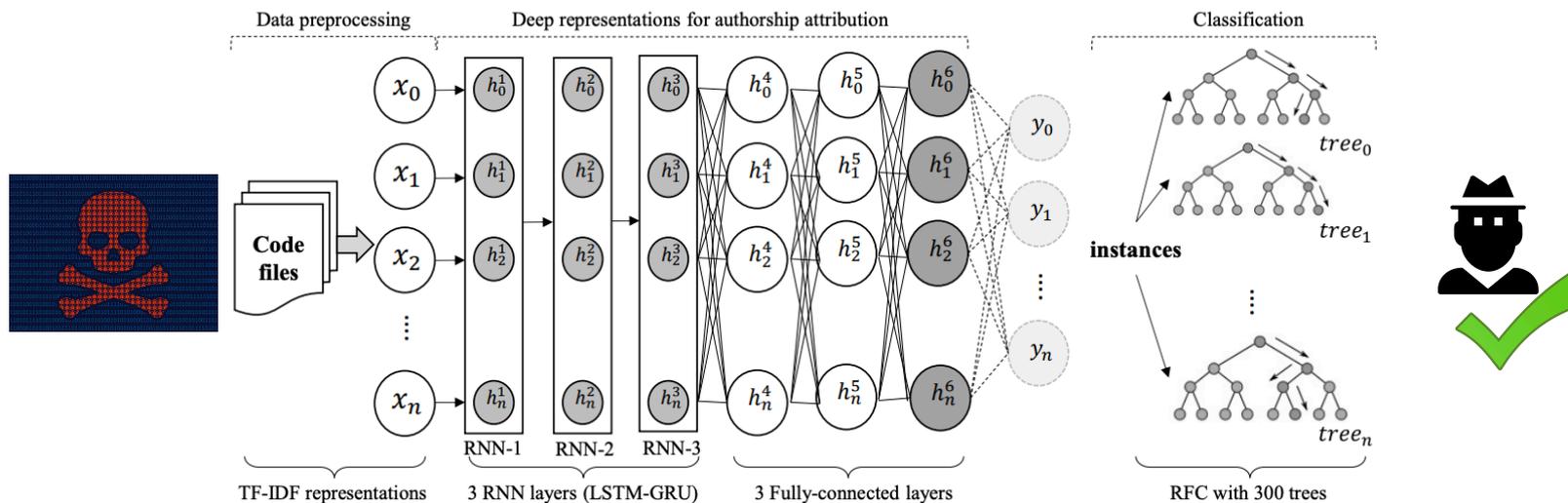
Putting ML for Good Use

- Security defense
 - Intrusion and anomaly detection
 - Deception detection
 - System failure detection
- Data synthesis
 - Generating “new malware that never exist before” → proactive defense
 - Vulnerability discovery, system patching
 - Designing new network protocols



Code Author Attribution

Large-Scale and Language-Oblivious Code Authorship Identification, CCS 2018

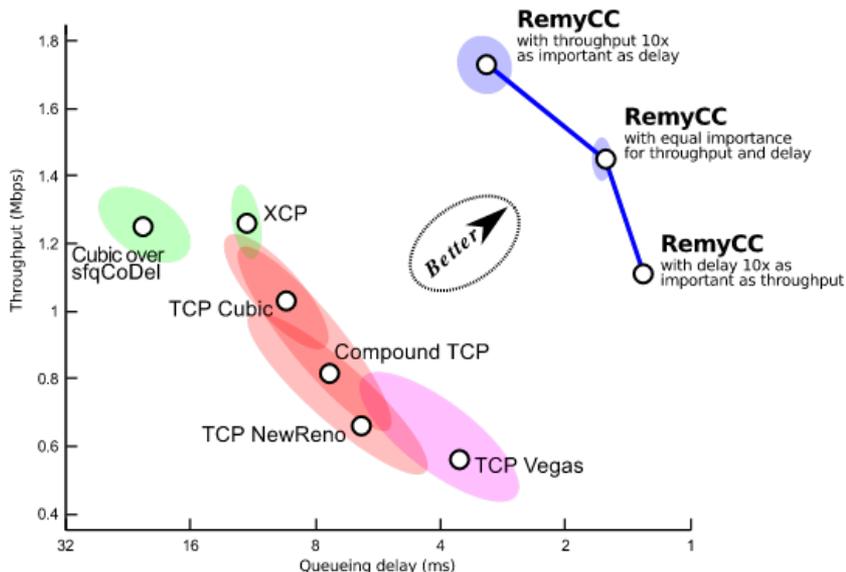


Another Example: TCP ex Machina

Computer-Generated Congestion Control, SIGCOMM'13

TCP congestion control basics:

- Slow start
- Congestion window
- Fast retransmit



Shown here: median results and 1- σ ellipses for eight endpoints contending for a 15 Mbps link, RTT = 150 ms, exponentially-distributed flow lengths and pause times.



Transparency and Biases

- How does ML made decisions?
- How robust are ML systems?

