tejasgokhale.com



# **Lecture 1: Introduction**

CMSC 491/691

Robust Machine Learning



tejasgokhale.com





# Lecture 1: Introduction

CMSC 491/691

Robust Machine Learning



### Welcome to the class!

This is an interactive class. Don't just consume – participate!

This is a research class. Don't just "learn" – critique, ideate, design, evaluate!

#### tejasgokhale.com

## Tejas Gokhale



Assistant Professor Computer Science University of Maryland, Baltimore County



2011—2015 B.E. (Honours) BITS Pilani (Goa) 2016-2018: M.S. Carnegie Mellon University





2023—present Assistant Professor University of Maryland Baltimore County



## Tejas Gokhale



Assistant Professor Computer Science University of Maryland, Baltimore County

### tejasgokhale.com TAY + JUSS

GO



2011—2015 B.E. (Honours) BITS Pilani (Goa) 2016-2018: M.S. Carnegie Mellon University

Temperature Mean: 10.4 °C Precipitation Sum: 970.1 mm

2018—2023: Ph.D. Arizona State University

**CLAY** 

Temperature Mean: 23.8 °C Precipitation Sum: 209.4 mm





2023—present Assistant Professor University of Maryland Baltimore County





#### Research Areas

- Computer Vision
- Vision & Language
- Image Generation
- Robustness & Reliability
- Active Perception

#### www.tejasgokhale.com



### Tejas Gokhale, Assistant Professor



### **Cognitive Vision Group, ITE 368**

#### **Current Projects**

- Domain Adaptation/Generalization
- Quantifying Visual Typicality
- Visual Semantics and Pragmatics
- Continual Concept Learning
- Generative Model Evaluation
- Multimodal Active Inference

#### Highlights / Activities

- AAAI 2024: New Faculty Highlights Invited Talk
- ECCV'24 & WACV'24: Tutorials on Reliability of Generative Models
- Area Chair: NeurIPS, ACL ARR, EMNLP, NAACL, WACV
- UMBC PPR Seminar









### **Course Staff**



Instructor: Tejas Gokhale Assistant Professor, CSEE OH: Wed 2:30 – 3:30 PM ITE 214 gokhale@umbc.edu TA/Grader ...

### Course Website

https://courses.cs.umbc.edu/graduate/691rml/



## **Quick Introductions**

#### (1) Name

- (2) Major (e.g. CS)
- (3) Level (BS/MS/PhD)
- (4) Why are you taking this class?
- (5) Something you did this summer





### **Course Website**



## **Class Structure: Overview**

- [Wed 08/28]: Today, Class Logistics and Introduction
- [Mon 09/02]: Labor Day, NO CLASS
- [Wed 09/04]: Machine Learning Review
- Every Week after that:
  - MON: [TEJAS] Overview of a Robustness Challenge
  - WED: [STUDENTS] Quiz; Paper Presentations; Class Discussion

### We have 12 Robustness Topics

## **Class Structure: Tentative Schedule**

### We have 12 Robustness Topics

	А	В	С	D	E	F	G	Н	I
1	Date	Day	Торіс	Notes		Date	Day	Торіс	Notes
2	26-Aug	Μ				28-Aug	W	Intro	
3	2-Sep	Μ	LABOR DAY			4-Sep	W	ML Refresher	
4	9-Sep	Μ	Domain Adaptation			11-Sep	W	Presentations	
5	16-Sep	Μ	Domain Generalization			18-Sep	W	Presentations	
6	23-Sep	Μ	OOD Detection			25-Sep	W	Presentations	
7	30-Sep	Μ	Adversarial Attacks, Backdoor Attacks			2-Oct	W	Presentations	
8	7-Oct	Μ	Uncertainty and Calibration			9-Oct	W	Presentations	
9	14-Oct	Μ	Online/Continual Learning			16-Oct	W	Presentations	
10	21-Oct	Μ	Self-Supervised/Unsupervised Learning			23-Oct	W	Presentations	
11	28-Oct	Μ	Test-Time Learning, Adaptation			30-Oct	W	Presentations	
12	4-Nov	Μ	Machine Unlearning, Model Editing			6-Nov	W	Presentations	
13	11-Nov	Μ	Interpretability, Explanability, Compositionality			13-Nov	W	Presentations	
14	18-Nov	Μ	ML Evaluation: Metrics a	3	20-Nov	W	Presentations		
15	25-Nov	Μ	Robustness Tradeoffs			27-Nov	W	no class (but take home as	ssignment?)
16	2-Dec	Μ	<b>Project Presentations</b>			4-Dec	W	<b>Project Presentations</b>	
17	9-Dec	Μ	Invited Talk / "The Last L						

## **Class Structure: Grading**



When the teacher says you won't be able to do all the homework in one night



## **Class Structure: Grading**

### We have 12 Robustness Topics

- Paper Presentations 20% pick two topics (10% each)
- Survey Papers 20% pick two topics ( $\neq$  your presentation topic)
- Quizzes (each Wed) 25% best 10 out of 12 (2.5% each)
- Project (group of 3)\* 35%
  - o Proposal (5%), Midterm Update Video (5%), Final Presentation (10%), Final Report (15%)
  - \*If you're a PhD/MS thesis student, you can opt to work alone
    - But you need my APPROVAL!

• *Extra Credit max 10%* (opportunities will be announced periodically)

## **Class Structure: Deadlines & Late Submission**

- **Paper Presentations:** In class (sign up sheet will be shared)
  - o If you are scheduled to present, but can't make it, send an email to me and the TA to request date-change
  - o Permissible reasons: unforeseen events (illness, injury, emergency), travel to academic conferences, interviews
  - Bottom Line: I reserve the right to approve or deny date-change requests
- Survey Papers: Due "next Wednesday 2359 UMBC time". See example below.
  - If Lecture on "Domain Adaptation" is on

- Monday, Sept 09.
- Then Survey paper on "Domain Adaptation" is due on Wednesday, Sept 18, 2359 UMBC time.
  Late submissions: 10% deducted for each late day.
- Quizzes: In class. 12 quizzes. Your highest 10 grades will be chosen.
  OBottom Line: if you miss a quiz, you miss a quiz.
- **Project:** Each milestone has fixed deadline.
  - $\circ$  Late Submissions: ~~ 10% deducted for each late day FOR ALL GROUP MEMBERS

## **Academic Integrity**

I take academic integrity very seriously. You should too. If you're unsure about something, ask us. You are at a top-tier (R1) research university. Use this privilege to learn. A good grade will follow. Don't throw away this opportunity by cheating.

- Presentations, Survey Papers, Quizzes must be done independently.
- Do not plagiarize. Consequences will not be pleasant.
- Do not use "AI" assistants for any part of any assignment. Consequences will not be pleasant.
- Familiarize yourself with UMBC policy on plagiarism and other forms of cheating: <u>https://academicconduct.umbc.edu/resources-for-students/</u>
- Read the syllabus for consequences of academic integrity violations.



## **Academic Integrity**

I take academic integrity very seriously. You should too. If you're unsure about something, ask us. You are at a top-tier (R1) research university. Use this privilege to learn. A good grade will follow. Don't throw away this opportunity by cheating.

#### 3.6 Good Practices

If the integrity of your work in this course is challenged, you are responsible for demonstrating proof that the work submitted is your own. A good starting point is to enable versioning/tracking in Google Docs, Word, Pages, or other software so that your writing activities/progress during the semester can be logged if necessary. Keeping copies of research notes, scribbles, and related material may be helpful, too.

#### 3.7 Viva or Oral Defense of Flagged Submissions

To ensure academic and professional integrity, I reserve the right to hold a one-on-one oral viva (defense) of submissions deemed questionable, to determine your knowledge and mastery of the topic/resources versus the material submitted. Failing that viva will result in an 'F' on the assignment and an Academic Integrity violation report filed with the Graduate School.

## **Academic Integrity**

I take academic integrity very seriously. You should too. If you're unsure about something, ask us. You are at a top-tier (R1) research university. Use this privilege to learn. A good grade will follow. Don't throw away this opportunity by cheating.

#### 3.8 Penalties

Academic misconduct could result in disciplinary action that may include, but is not limited to, suspension or dismissal. The **absolute minimum penalty** for a first offense of academic dishonesty in this course is a grade of zero on the assignment and a one-letter-grade reduction in the final class grade. However, depending on the nature of the offense, the penalty may be more severe, including but not limited to an **F** for the course, suspension, or expulsion. The minimum penalty for a second offense of academic dishonesty is an **F** for the course without possibility of dropping, but may be more severe.

## I DON'T ALWAYS CARE About My Grade...

### BUT WHEN I DO IT'S THE END OF THE SEMESTER AND EVEN THOUGH I DIDN'T DO ALL THE ASSIGNMENTS, I'LL ASK FOR EXTRA CREDIT NOW.

## Seek Help Early!



### Help us help you.

## 491 (Undergrad) vs 691 (Grad)

This is a "Research" class – anyone with the right mindset (+ an intro ML class) is ready

#### Main difference: projects (scope and novelty)

- Projects will be graded in terms of "relative growth"
  - Some may have previous research experience
  - **o** Some are taking this class to get research experience
- Grad projects:
  - $_{\odot}$  Original & unique research hypothesis with a potential for publication
- Undergrads projects can be:
  - $\circ$  Original & unique research hypothesis with a potential for publication
  - $_{\odot}$  Working on an idea that we provide (i.e. you get to skip ''ideation'')
  - $\,\circ\,$  Innovative applications or combination of existing work

Grad student and undergrad sitting in the same class



## FAQ: Can I join your research lab?

• Joining (See FAQ on my website)

 $\odot$  Take this class and talk to me during office hours about your interests

- Will I get paid for research ?
  - Depends (I currently only have funding for PhD students)
  - o Undergraduate Research & Prestigious Scholarships UMBC
  - o <u>CWIT Scholars Center for Women in Technology UMBC</u>
  - You can also do research for credit (e.g. CMSC 299, 499, 698, 699)

## **Other Questions?**





Robust Machine Learning Lecture 0: Introduction





### Models that learn from data are embedded in our lives













### Models that learn from data are embedded in our lives

Recent advances have been rapidly adopted by common, non-expert users

#### DALL·E Now Available Without Waitlist

New users can start creating straight away. Lessons learned from deployment and improvements to our safety systems make wider availability possible.

SIGN UP 7



THE SHIFT

#### An A.I.-Generated Picture Won an Art Prize. Artists Aren't Happy.

"I won, and I didn't break any rules," the artwork's creator says.

#### Forbes

FORBES > LEADERSHIP > CAREERS

#### Educators Battle Plagiarism As 89% Of Students Admit To Using OpenAI's ChatGPT For Homework



#### MICROSOFT / TECH / ARTIFICIAL INTELLIGENCE

#### Microsoft announces new Bing and Edge browser powered by upgraded ChatGPT AI

/ Microsoft says it's using conversational AI to create a new way to browse the web. Users will be able to chat to Bing like ChatGPT, asking questions and receiving answers in natural language.

## ImageNet: An ML Home Run





But what do these results *really* mean?

## A Limitation of the (Supervised) ML Framework



#### Measure of performance:

Fraction of mistakes during testing

**But:** In reality, the distributions we **use** ML on are NOT the ones we **train** it on

## A Limitation of the (Supervised) ML Framework



**Measure of performance:** Fraction of mistakes during testing

**But:** In reality, the distributions we **use** ML on are NOT the ones we **train** it on

What can go wrong?

### Standard i.i.d. Assumption in Machine Learning



"Independent and Identically Distributed" Models learn useful patterns

### Standard i.i.d. Assumption in Machine Learning



#### IID Assumption collapses in real-world "in-the-wild" settings Model performance deteriorates

# **Robustness in Computer Vision**

Findings from Previous Work

## Poses can fool Image Classifiers



school bus 1.0 garbage truck 0.99 punching bag 1.0 snowplow 0.92

Alcorn et al. CVPR 2019

# Poses can fool Image Classifiers

- Goal: correctly classify previously unseen test images.
- Statistical ML operates with the "i.i.d." assumption
- But real-world test inputs are often NOT i.i.d. !!!
  - Poses can fool classifiers
    - Rotation
    - Translation
    - Scale
    - Occlusion
    - ...



school bus 1.0 garbage truck 0.99 punching bag 1.0 snowplow 0.92

Alcorn et al. CVPR 2019

## Natural Corruptions affect accuracy



Hendrycks et al. ICLR 2019

## **Spurious Correlations / Biased Datasets**

#### **Common training examples** Test examples y: landbird y: waterbird y: waterbird a: land a: land a: water background background background Waterbirds y: dark hair y: blond hair ·UCUS y: blond hair LUUNUING SPU a: female a: male a: male CelebA

Sagawa et al. ICLR 2020
# **Adversarial Attacks on Image Classifiers**

- Algorithms that can "find" perturbations to add to images, in order to fool classifiers
- Given image x, find g(x) s.t.  $x + \epsilon g(x)$  fools classifier
- Perturbations are typically norm-bounded



Goodfellow et al. ICLR 2015

#### Lack of Diverse Data hurts Reliability

Mallard Water Background

Penguin



Mallard Snow Background

#### Lack of Diverse Data hurts Reliability

Mallard Snow Background





# **Domain Shift is a Nuisance**



Single-Source Domain Generalization for Digit Classification



# ML Predictions Are (Mostly) Accurate but Brittle



[Szegedy Zaremba Sutskever Bruna Erhan Goodfellow Fergus 2013] [Biggio Corona Maiorca Nelson Srndic Laskov Giacinto Roli 2013]

**But also:** [Dalvi Domingos Mausam Sanghai Verma 2004][Lowd Meek 2005] [Globerson Roweis 2006][Kolcz Teo 2009][Barreno Nelson Rubinstein Joseph Tygar 2010] [Biggio Fumera Roli 2010][Biggio Fumera Roli 2014][Srndic Laskov 2013]

# Biker Pedestrian Sign





#### Persons









Green Traffic Light

Adversarial Perturbation Attack

#### Pedestrian Sign





Speed Limit 45 Sign

# Adversarial Rotation Attack









Adversarial Patch Attack

# **Adversarial Examples**

• In 2014, one of the seminal papers of Goodfellow et al. shows that an adversarial image of a panda can fool the ML model to output "gibbon", which started the area of adversarial ML

Original image



Classified as panda 57.7% confidence

Adversarial image





Gibbon

Small adversarial noise

Classified as gibbon 99.3% confidence

# **Adversarial Examples**

• Similar example, from Szagedy et al. (2014)



Picture from: Szagedy et al. (2014) – Intriguing Properties of Neural Networks





# **Adversarial Attacks**

Algorithms that can "find" perturbations to add to images, in order to fool classifiers

Given image x, find g(x) s.t.  $x + \epsilon g(x)$  fools classifier Perturbations are typically norm-bounded



Goodfellow et al. ICLR 2015

# **Adversarial Training**

Leverages the concept of adversarial examples, in order to improve classifier robustness to such attacks

#### min—max optimization

maximization: find adversarial images minimization: train classifier to correctly classify such images norm-bounded perturbations ==> robustness within the norm-ball



$$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \left[ \max_{\delta\in\mathcal{S}} L(\theta, x + \delta, y) \right]$$

Madry et al. ICLR 2018

# Physical-World Attack: Printed Adversarial Images

 Not only adversarial examples in the digital world, but printed adversarial images can also fool machine learning models



# Physical-World Attack: Adversarial STOP Sign

- An example of manipulating a STOP sign with adversarial patches
  - Methodology: carefully design a patch and attach it to the STOP sign
  - Cause the DL model of a self-driving car to misclassify it as a Speed Limit 45 sign
    - The authors achieved 100% attack success in lab test, and 85% in field test

#### Lab (Stationary) Test

Physical road signs with adversarial perturbation under different conditions





Stop Sign → Speed Limit Sign

#### Field (Drive-By) Test

Video sequences taken under different driving speeds





Stop Sign → Speed Limit Sign

Picture from: Eykholt (2017) - Robust Physical-World Attacks on Deep Learning Visual Classification

# Physical-World Attack: Adversarial STOP Sign

More examples of lab test for STOP signs with a target class Speed Limit 45



Picture from: Eykholt (2017) - Robust Physical-World Attacks on Deep Learning Visual Classification

# Physical-World Attack: Adversarial Patch

- Not only adversarial patch can fool a classifier, but also a SOTA detector
- An example of a person wearing an adversarial patch who cannot be detected by a YOLOv2 model
  - This can be used by intruders to get past security cameras



Thys et al. (2019) - Fooling automated surveillance cameras: adversarial patches to attack person detection

#### Physical-World Attack: Attack Tesla Autopilot System

 Non-scientific example: a Tesla owner checks if the car can distinguish a person wearing a cover-up from a traffic cone



# Why should we care?

 $\rightarrow$  People suffer consequences because of use in real-world systems  $\rightarrow$  Safety, security, trust in the systems that we engineer



#### Wrongfully Accused by an Algorithm

In what may be the first known case of its kind, a faulty facial recognition match led to a Michigan man's arrest for a crime he did not commit.





TECHNOLOGY

Feds Say Self-Driving Uber SUV Did Not Recognize Jaywalking Pedestrian In Fatal Crash

# Data Poisoning

#### **Goal:** Maintain training accuracy but hamper generalization



# Data Poisoning

#### Goal: Maintain training accuracy but hamper generalization



- → Fundamental problem in "classic" ML (robust statistics)
- → But: seems less so in deep learning
- → **Reason:** Memorization?

# Data Poisoning

#### classification of **specific** inputs

### Goal: Maintain training accuracy but hamper generalization



"van"

"dog"

[Koh Liang 2017]: Can manipulate many predictions with a single "poisoned" input

#### But: This gets (much) worse

[Gu Dolan-Gavitt Garg 2017][Turner Tsipras M 2018]: Can plant an **undetectable backdoor** that gives an almost **total** control over the model

(To learn more about backdoor attacks: See poster #148 on Wed [Tran Li M 2018])

# You don't even need poisonous samples.

## Re-ordering training batches of clean data $\rightarrow$ failures





Figure 1. Using an external model such as CLIP, the distribution of the confounding class (truck) samples is examined with respect to their softmax probabilities of the attacked class (car). The samples with the highest probability are used to train an image classifier, which confounds the classifier's ability to distinguish between the attacked and confounding classes.

# Are we doomed?



# Are we doomed?

(Is ML inherently not reliable?)

# Are we doomed? (Is ML inherently not reliable?)

We need to re-think how we do ML





OCTOBER 30, 2023

### Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence

BRIEFING ROOM > PRESIDENTIAL ACTIONS

By the authority vested in me as President by the Constitution and the laws of the United States of America, it is hereby ordered as follows:



#### **Robustness Definitions**

(or the lack thereof)



# Adversarial vs. Out-of-distribution

• Each of the described attacks can further be:

Adversarial example (algorithmically manipulated example)



Small adversarial noise

Classified as panda 57.7% confidence

Adversarial image



Classified as gibbon 99.3% confidence

• **Out-of-distribution** example (natural example)





# Robustness

Predictor ff(x) = y for  $x \in \mathcal{X}$ 

#### **Robustness to small perturbations:** $f(x + \Delta x) = ?$ s.t. $||\Delta x||_p \le \epsilon$



#### **Robustness to small transformations:**

$$(g(x)) = ?$$
 s.t.  $||g(x) - x||_p \le \rho$   
and  $g: \mathcal{X} \to \mathcal{X}$ 



Examples from "ImageNet-C" and "ImageNet-P" datasets – Hendrycks et al. ICLR 2019

# Generalization

#### **IID Generalization**

Training and Testing inputs sampled independently and identically from a common underlying distribution  ${\mathcal X}$ 

$$\begin{array}{l} x_{train} \sim \mathcal{X} \\ x_{test} \sim \mathcal{X} \end{array}$$

#### OOD Generalization aka Domain Generalization

Training and Testing inputs samples from non-identical distributions  $\mathcal{X}_{source}$  and  $\mathcal{X}_{target}$ 

Thus test images are "out-of-distribution" for the classifier that has seen training images.

$$\begin{array}{l} x_{train} \sim \mathcal{X}_{source} \\ x_{test} \sim \mathcal{X}_{target} \end{array}$$

#### Some of Tejas' Research

#### (or why I care about Robust Machine Learning)

# **Perception & Reasoning with Robustness**

#### **Robust Image Recognition**



Effects of multiple data sources on OOD and adversarial robustness

Gokhale AAAI'21; Gokhale ACL'22; Gokhale WACV'23; Cheng ICCV '23; Wisdom arxiv 2023; Kulkarni CVPR-W'21



#### Scene Completion for Missing Sensor/Modality



# **Perception & Reasoning with Robustness**

#### Robust Visual Reasoning (Visual QA, Video Captioning, V&L Inference)

V&L Robustness: Logical, Semantic, Spatial (use additional knowledge sources and sensors)









Gokhale ECCV '20; Gokhale EMNLP'20; Gokhale ACL'21; Fang EMNLP'20; Banerjee ICCV'21; Patel EMNLP'22

# **Perception & Reasoning with Robustness**

#### Robust Visual Reasoning (Visual QA, Video Captioning, V&L Inference)

V&L Robustness: Logical, Semantic, Spatial (use additional knowledge sources and sensors)



Is the for	k <b>NOT</b> on the plate?	*10
yes	94.785%	Negation
no	5.215%	
s the for	k on the plate AND is the food	made of eggs
5 the for	con the place, true is the lood	induc of eggs
no	97.855%	Conjunct
yes	a 144%	
the fork	on the plate <b>OR</b> is the food ma	de of eggs?
yes <b>s the fork</b> 10	e on the plate OR is the food ma	de of eggs?



Answer s that a giraffe or an elephant? Giraffe Who is feeding the giraffe *behind* the Ladv s there a fence *near* the animal *behind* Yes On which side of the image is the man? s the giraffe behind the man?





Understanding Agent Actions in Videos with Commonsense, Counterfactual and Physics-Based Reasoning



Counterfactual Question What will happen if the yellow cube is removed ?

(A) Purple Cube will collide with brown cube

Planning Question How can the collision between yellow and purple cube be stopped?





**Conventional Caption** Group of runners get prepared to run a race. Commonsense-Enriched In order to win a medal, a group of runners get prepared to run a race. As a

Caption

result they are congratulated at the finish line. They are athletic.

**Commonsense Question** Answering

What happens next to the runners?

Are congratulated at the finish line become tired

Gokhale ECCV '20; Gokhale EMNLP'20; Gokhale ACL'21; Fang EMNLP'20; Banerjee ICCV'21; Patel EMNLP'22
# **Concept Learning/Distillation in Text-to-Image Generation**

#### Few-Shot Concept Leaning in Text-to-Image Models

- learn common visual concepts from a few images •
- assign semantic meaning (in latent space) •
- Reproduce the concept (novel view synthesis)



Patel AAAI '24; Gokhale Tech Rep 2022; Chatterjee ECCV '24; Chatterjee ECCV '24

Spatial Reasoning in Text-to-Image Models

Evaluation Dataset (SR2D) and Metrics (VISOR)







Generated Image x=g(t)

- A = motorcycle, B = elephantLocate Centroids of A, B
- Improving T2I spatial reasoning with
  - 1. Recaptioning

(how can a very small number of specialized "spatial" captions" help finetune T2I models?)

### 2. guidance from graphics engines



# **Concept Learning/Distillation in Text-to-Image Generation**

Few-Shot Concept Leaning in Text-to-Image Models

- learn common visual concepts from a few images
- assign semantic meaning (in latent space)
- Reproduce the concept (novel view synthesis)



Patel AAAI '24; Gokhale Tech Rep 2022; Chatterjee ECCV '24; Chatterjee ECCV '24

#### Spatial Reasoning in Text-to-Image Models

• Evaluation Dataset (SR2D) and Metrics (VISOR)



- Improving T2I spatial reasoning with
  - 1. Recaptioning

(how can a very small number of specialized "spatial" captions" help finetune T2I models?)

### 2. guidance from graphics engines



## Next Time:

## Machine Learning Review

- Machine Learning Training and Inference Pipeline
- Neural Networks / Deep Learning
- ML methods for visual recognition