

Lecture Notes for **Machine Learning in Python**

Professor Eric Larson
Final Lecture: Ethics and Retrospective

Lecture Agenda

- Logistics
 - Grading Update
 - Sequential Networks due **Last Day of Finals**
 - **Before NOON on December 13**
- Agenda
 - Ethical Principles
 - Retrospective and Evaluations

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

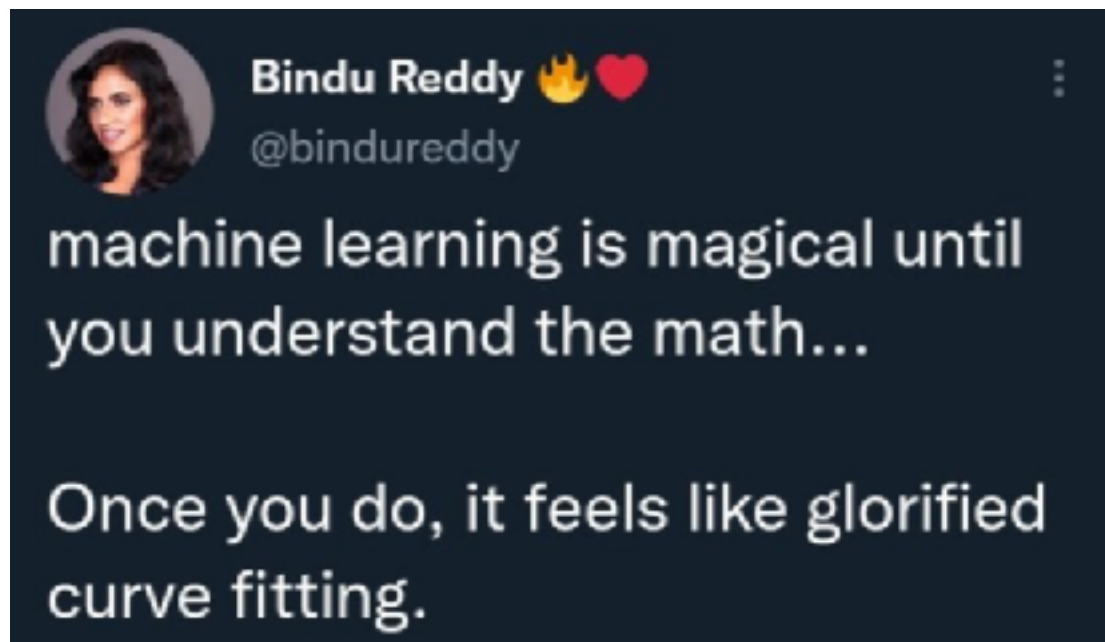
Sequential
Networks

Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

Sequential Networks Town Hall



AI Ethics Principles



Janelle Shane @JanelleCShane · 1d
Predictive policing algorithms don't predict who commits crime. They predict who the police will arrest.



Emily M. Bender, professionally... · 11h ...

"AI" can NOT:
* Predict who will commit a crime

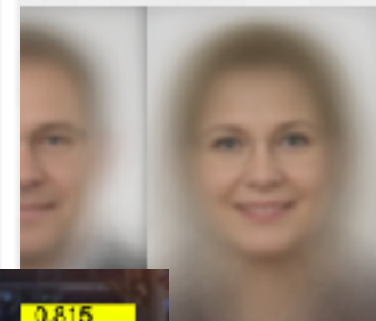
"AI" can:
* Make biased policing look "objective"



Timnit Gebru: Gender Shades



Lighter Female	Largest Gap
98.2%	20.8%
94.0%	33.8%
92.9%	34.4%



Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*

Joy Buolamwini
MIT Media Lab 75 Amherst St. Cambridge, MA 02139

JOYBU@MIT.EDU

Timnit Gebru
Microsoft Research 641 Avenue of the Americas, New York, NY 10011

TIMNIT.GEBRU@MICROSOFT.COM

<http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf>



- **Reliability:** reliably operate in accordance with their intended purpose
- **Beneficence:** individuals, society and the environment.
- **Respect:** respect human rights, diversity, and autonomy of individuals.
- **Fairness:** be inclusive and accessible, and should not involve or result in unfair discrimination against individuals, communities or groups
- **Privacy:** respect and uphold privacy rights and data protection, and ensure the security of data
- **Transparency:** ensure people know when they are being significantly impacted by an AI system, and can find out when engaging with them
- **Contestability:** should be a timely process to allow people to challenge the use or output of the AI system
- **Accountability:** Those responsible for the different phases of the AI system lifecycle should be identifiable and accountable for the outcomes of the AI systems, and *human oversight* of AI systems should be enabled.

To enforce these principles a board with autonomy must exist

Bias Case Study in NLP



Timnit Gebru ✓
@timnitGebru

I'm sick of this framing. Tired of it. Many people have tried to explain, many scholars. Listen to us. You can't just reduce harms caused by ML to dataset bias.



Yann LeCun @ylecun · 19h

ML systems are biased when data is biased. This face upsampling system makes everyone look white because the network was pretrained on FlickrFaceHQ, which mainly contains white people pics....

✓ **Dataset Bias:** Over-representing a specific group of data, potentially leading to performance differences across groups.

ML Fairness: Understanding and considering the harms that performance differences can incur on a specific group.

Example:

- A facial identification system used by police has a 1.2% error rate.
- For white individuals this error is 0.8%
- For black individuals this error is 1.9%
- The models are retrained across groups and now the error rate is 1.4% across all groups.
- Is the system fair?



François Chollet ✓ @fchollet · 11h

When faced with tech ethics problems, you can either ask hard questions, seek solutions, and take responsibility, or you



Devin Guillory @databoydg · 13h

Watching one of the most influential

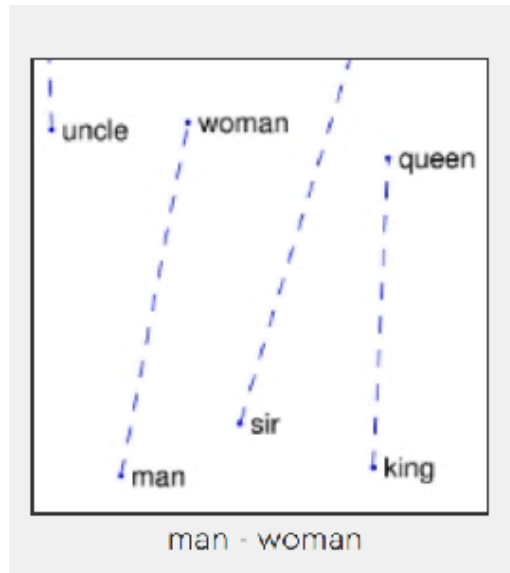
Timnit Gebru

A lot of times, people are talking about bias in the sense of equalizing performance across groups. They're not thinking about the underlying foundation, whether a task should exist in the first place, who creates it, who will deploy it on which population, who owns the data, and how is it used?

The root of these problems is not only technological. It's social.

Using technology with this underlying social foundation often advances the worst possible things that are happening. In order for technology not to do that, you have to work on the underlying foundation as well. You can't just close your eyes and say: "Oh, whatever, the foundation, I'm a scientist. All I'm going to do is math."

Ethics in Language: Gender Inequity



Trained on
New York Times



$$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"aunt"}) - W(\text{"uncle"})$$

$$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"queen"}) - W(\text{"king"})$$

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{computer programmer}} - \overrightarrow{\text{homemaker}}$$

Extreme *she* occupations

- | | | |
|-----------------|-----------------------|------------------------|
| 1. homemaker | 2. nurse | 3. receptionist |
| 4. librarian | 5. socialite | 6. hairdresser |
| 7. nanny | 8. bookkeeper | 9. stylist |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

Extreme *he* occupations

- | | | |
|----------------|-------------------|----------------|
| 1. maestro | 2. skipper | 3. protege |
| 4. philosopher | 5. captain | 6. architect |
| 7. financier | 8. warrior | 9. broadcaster |
| 10. magician | 11. fighter pilot | 12. boss |

Bolukbasi et al., NeurIPS 2016

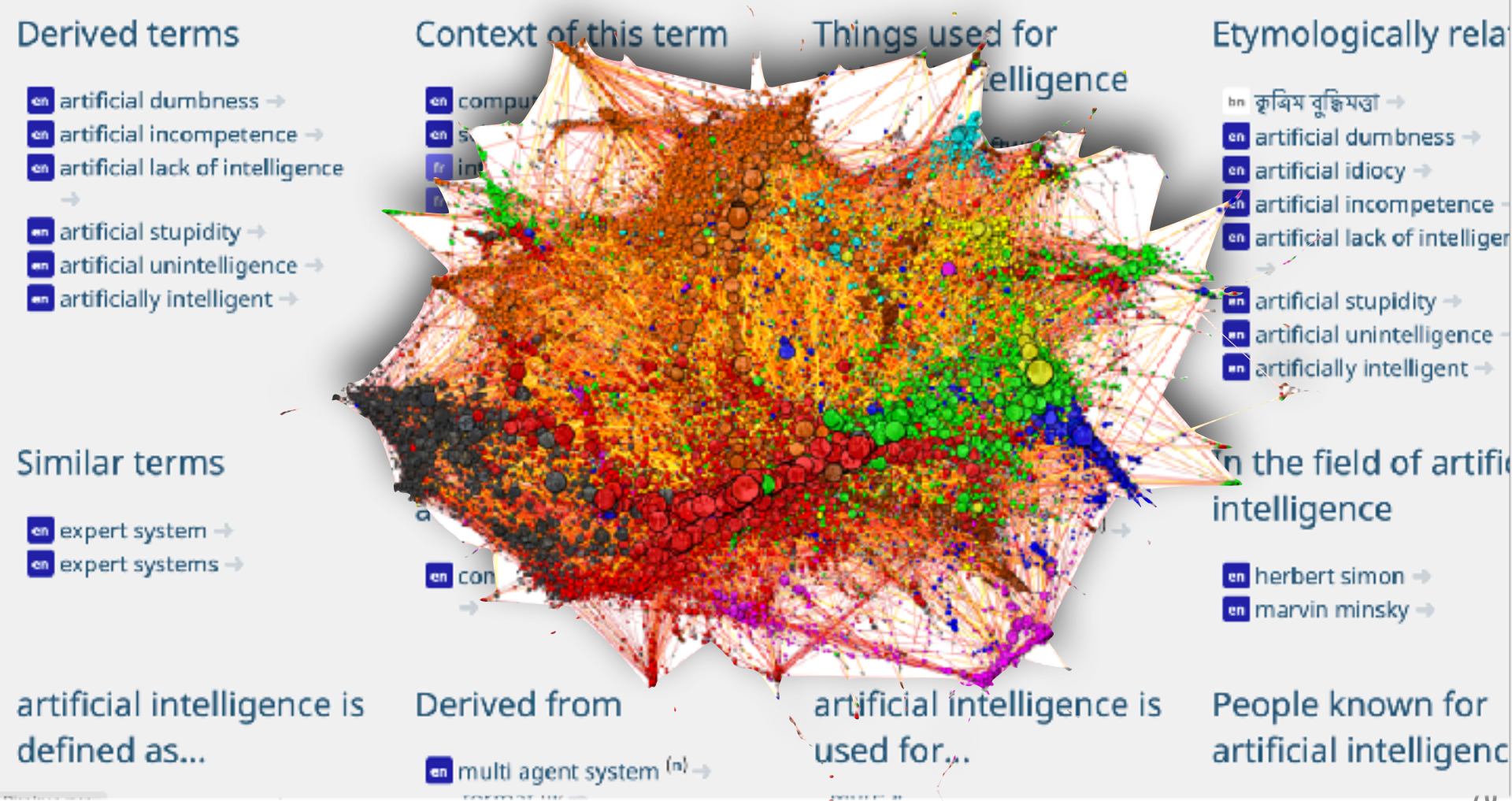
<https://arxiv.org/pdf/1607.06520.pdf>

<https://nlp.stanford.edu/projects/glove/>

ConceptNet, a Multi-lingual Knowledge Graph

artificial intelligence

An English term in ConceptNet 5.8



ConceptNet Numberbatch



- **Step One:** Create a Knowledge Graph (from multiple sources with relations like *UsedFor*, *PartOf*, etc.)
- **Step Two:** Based on this KG, perturb existing embeddings (like GloVe) to minimize:

$$\Psi(Q) = \sum_{i=1}^n \left[\underbrace{\alpha_i \| \underset{\substack{\uparrow \\ \text{new embed}}}{q_i} - \underset{\substack{\uparrow \\ \text{old embed}}}{\hat{q}_i} \|^2}_{\text{(keep similar to original)}} + \underbrace{\sum_{(i,j) \in E} \beta_{ij} \| q_i - q_j \|^2}_{\substack{\text{neighbors from KG} \\ \text{(make similar according to other knowledge)}}} \right]$$

- Straight forward to optimize the objective by averaging neighbors in the ConceptNet Knowledge Graph
- Multiple embeddings achieved by merging through “retrofitting” which projects onto a shared matrix space (with SVD)

Lightning Demo (or self guided demo)



How to Make a Racist AI without Really Trying



Robyn Speer, 2017

<http://blog.conceptnet.io/posts/2017/how-to-make-a-racist-ai-without-really-trying/>

Debiasing: Man is to Computer Programmer as Woman is to Homemaker? De-biasing Word Embeddings

Bolukbasi et al., NeurIPS 2016

<https://arxiv.org/pdf/1607.06520.pdf>

ConceptNet 5.5: An Open Multilingual Graph of General Knowledge

Speer et al., AAAI 2017

<https://arxiv.org/pdf/1612.03975.pdf>

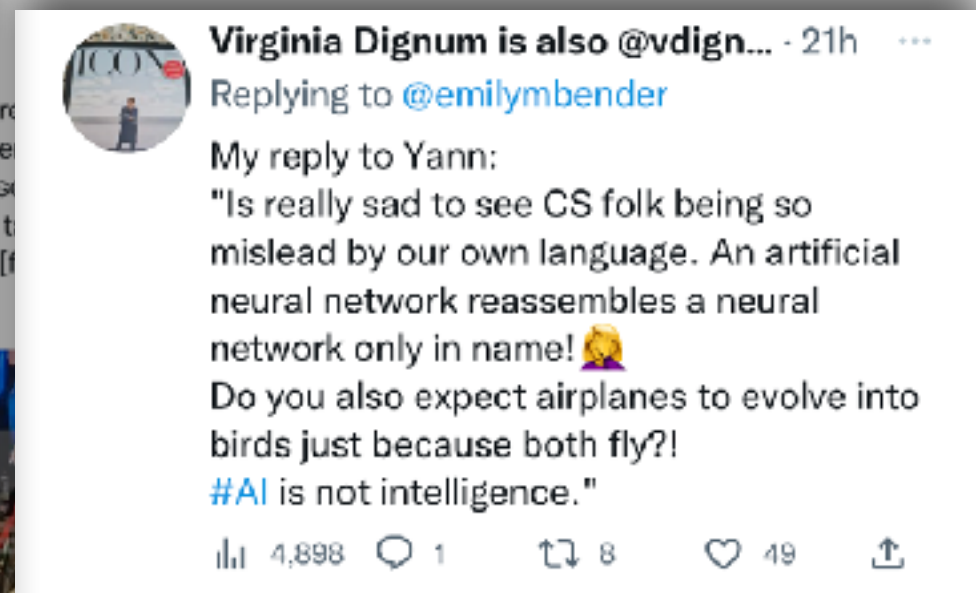


Rachael Tatman @rctatman · 18h

I first got interested in ethics in NLP/ML because I was asking "does this system work well for everyone". It's a good question, but there's a more important one:

Who is being harmed and who is benefiting from this system existing in the first place?

Course Retrospective



Course Retrospective

- AI winters exist
- machine learning

and history

- Formal methods
- At the end of

- **Open source** advancements

- <http://www>

Leading ML researchers issue statement of support for JMLR

From: Michael Jordan [mailto:jordan@CS.Berkeley.EDU]
Sent: Monday, October 08, 2001 5:33 PM
Subject: letter of resignation from Machine Learning journal

Dear colleagues in machine learning,

The forty people whose names appear below have resigned from the Editorial Board of the Machine Learning Journal (MLJ). We would like to make our resignations public, to explain the rationale for our action, and to indicate some of the implications that we see for members of the machine learning community worldwide.

The machine learning community has come of age during a period of enormous change in the way that research publications are circulated. Fifteen years ago research papers did not circulate easily, and as with other research communities we were fortunate that a viable commercial publishing model was in place so that the fledgling MLJ could begin to circulate. The needs of the community, principally those of seeing our published papers circulate as widely and rapidly as possible, and the business model of commercial publishers were in harmony.

Times have changed. Articles now circulate easily via the Internet, but unfortunately MLJ publications are under restricted access. Universities and research centers can pay a yearly fee of \$1050 US to obtain unrestricted access to MLJ articles (and individuals can pay \$120 US). While these fees provide access for institutions and individuals who can afford them, we feel that they also have the effect of limiting contact between the current machine learning community and the potentially much larger community of researchers worldwide whose participation in our field should be the fruit of the modern Internet.

None of the revenue stream from the journal makes its way back to authors, and in this context authors should expect a particularly favorable return on their intellectual contribution---they should expect a service that maximizes the distribution of their work. We see little benefit accruing to our community from a mechanism that ensures revenue for a third party by restricting the communication channel between authors and readers.

Sincerely yours,

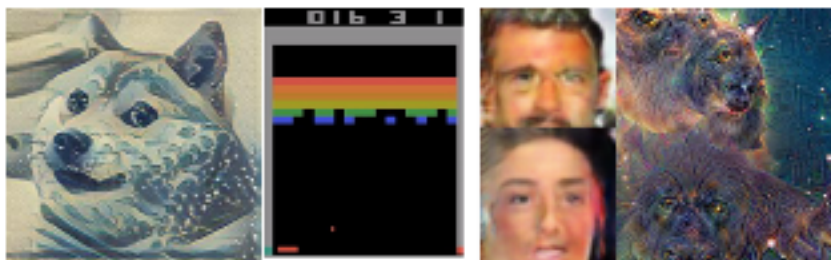
Chris Atkeson
Peter Bartlett
Andrew Barto
Jonathan Baxter
Yoshua Bengio
Kristin Bennett
Chris Bishop
Justin Boyan
Carla Brodley
Claire Cardie
William Cohen
Peter Dayan
Tom Dietterich
Jerome Friedman
Nir Friedman
Zoubin Ghahramani
David Heckerman
Geoffrey Hinton
Haym Hirsh
Tommi Jaakkola
Michael Jordan
Leslie Kaelbling
Daphne Koller
John Lafferty
Bridhar Mahadevan
Marina Meila
Andrew McCallum
Tom Mitchell
Stuart Russell
Lawrence Saul
Bernhard Schölkopf
John Shawe-Taylor
Yoram Singer
Satinder Singh
Padhraic Smyth
Richard Sutton
Sebastian Thrun
Manfred Warmuth
Chris Williams
Robert Williamson

Topics review

- Data **munging** in pandas and numpy and **visualization** with matplotlib, pandas, seaborn
- Data preprocessing: **dim reduction**, images, text, categorical features, **embeddings**
- **Linear models**: linear regression, logistic regression, simple neural networks
- **Optimization** strategies: Gradient ascent, Quasi-Newton, Extensions of SGD (RMSProp, AdaM)
- **Back propagation** in MLP (from scratch)
- Tensorflow/Keras for **wide and deep networks**
- **Convolutional** neural networks (up to modern day)
- **Sequential** neural networks (scratched surface only)

Topics Not Covered

- Transfer/Multi-Task Learning
- Visualizing Deep Convolutional Networks
- Fully Convolutional Networks
- Style Transfer (if time)
- Generative Networks
- Large Language Models



Syllabus for CSE8321: Machine Learning and Neural Networks

Course Schedule

Week	Lecture A	Lecture B	Lecture C
1	Lecture: Course Introduction and Syllabus	Lecture: Basics of Neural Networks	
2	Student Presentation and Reading: Deep Learning: Global 2017	Lecture: CNN Architecture Overview	
3	Lecture: CNN Visuals	Lecture: Image Style Transfer Overview	Lecture: CNN Visuals
4	Student Presentation and Reading: A Tale of Algorithms in AI (Part 1) (Fall 2018)	Student Presentation and Reading: Transfer and Style Transfer, 2018	
5	Lecture: Image Style Transfer in Detail	Lecture: Transfer Learning in CNNs	Lecture: Style Transfer
6	Lecture: Image Style Transfer in Detail	Lecture: Multi-modal Learning Overview	
7	Student Presentation and Reading: Deep Multitask Learning for Recommendation and Search, 2017	Student Presentation and Reading: An Overview of GANs: From Learning to Education (Fall 2018)	
8	Lecture: Generative Adversarial Networks	Lecture: Generative Adversarial Networks	Lecture: GANs
9	Student Presentation and Reading: Recent Progress in Deep Learning: A Survey	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
10	Lecture: Deep Reinforcement Learning Overview	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
11	Lecture: Deep Reinforcement Learning Overview	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
12	Student Presentation and Reading: Deep Reinforcement Learning Overview	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
13	Lecture: The Future of Deep Learning	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
14	Student Presentation and Reading: Deep Reinforcement Learning Overview	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs
15	Student Presentation and Reading: Deep Reinforcement Learning Overview	Lecture: GANs: From Learning to Education (Fall 2018)	Lecture: GANs

Syllabus for CSE8321: Machine Learning and Neural Networks

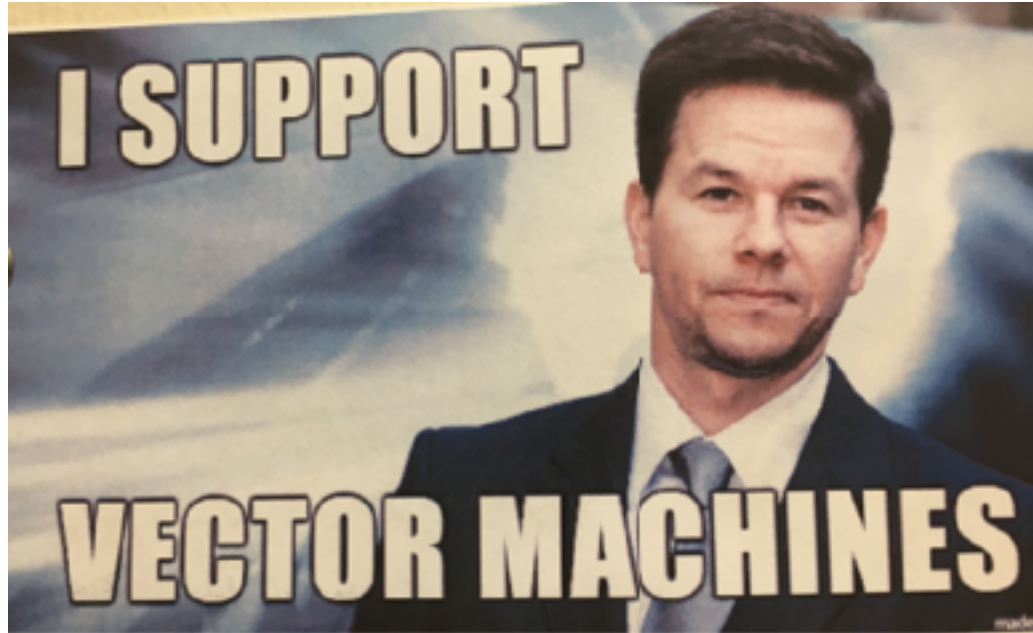
Overview

This course extends basic knowledge of the use of Neural Networks in machine learning beyonds simple prediction, especially targeted outputs that are generation or alteration of images, text, and audio. This course emphasizes topics of neural networks in the "deep learning" subdomain. This course will survey of important topics and current areas of research, including transfer learning, multi-task and multi-modal learning, image style transfer, neural network visualization, deep convolutional generative adversarial networks, and deep reinforcement learning. For grading, students are expected to complete smaller team-based projects throughout the semester, present one research paper in a 15-20 minute group presentation (covering topics in the course), and complete a comprehensive final project that involves a number of different deep learning architectures.

Thank you for a great semester!

- but it could **have been better** somehow, right?
 - how could you learn better, more reliably for an interview?
 - what should **not be cut** or **not changed**?
 - **Already cut:** SVMs, Ensembles, RNNs, many-to-many RNNs,
 - How Did X-formers go?
 - More convolutional approaches/depth?
 - More APIs? Turi / PyTorch?
 - More flipped Assignments?
 - Self-guided Jupyter notebooks?

Thank You for an Excellent Semester!



Courtesy of Omar Roa

Please fill out the course evaluations!!!!