Embeddings Projections / alignment

Adversarial Examples Adversarial Feature Learning

Ethics

Yoav Goldberg





About the exam





- 3 hours (maybe not all of it will be needed)
- Open materials

everything said in class is fair game. and also elaborations of the ideas in class.



- You need to be able to read and write ML models as precise mathematical formulas.
- You need to be able to suggest modification to models, and express them in writing / math.
- Calculate sizes of networks, rough estimate on runtime (which operations are expensive which are cheap), number of parameters of different models.





- "Will this network run faster on a CPU or a GPU?"
- Which changes to this network will make it work faster without sacrificing much power?



- You need to be able to give names to equation.
 - "this is a differentiable gate".
 - "this is an MLP".

. . .

- "this is attention".
- "this is regularization".

- You need to be able to explain concepts ("what is dropout?")
- "here are the equations of an LSTM. Modify them to add regularization on the bias term of the input gate".





What can and cannot be optimized with gradient descent?





- Think about loss functions
- Think about what a given architecture or representation can and cannot do
- "which of these systems produced this output"?





 Be able to think critically about design choices and their implications.



- Window vs RNN vs CNN vs Dilated CNN vs...
- Pre-trained embeddings, how to use them?
- How to represent words?
- LSTM Transducer vs Seq2Seq
- Attention vs. no Attention
- Softmax+cross-entropy vs margin loss
- Different pre-training architectures / choices



• "someone suggested to take architecture X and replace equation number (4) with

is it a good/bad idea? why? what will happen? do you need to do additional changes for this to work?"



- What do learning curves mean?
- How would you change a system to perform better? to do X?
- "Here is a system with some parameters. Here is a learning curve. Here is another learning curve.
 Which parameters did we change to get to the new curve?"

Embedding Projections



- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





- "Exploiting Similarities among Languages for Machine Translation" - Mikolov, Le & Sutskever, 2013
- Observed a similar structure in unsupervised embedding spaces of different languages, after rotation
- Learned a rotation matrix to translate words from one embedding space to another with some success
- Weakly supervised requires a small dictionary (5000 entries)





• Learned a rotation matrix to translate words from one embedding space to another with some success

$$\arg\min_{M} \sum_{x_i, y_i \in pairs} ||Mx_i - y_i||_2^2$$



• Learned a rotation matrix to translate words from one embedding space to another with some success

$$\arg \min_{M} \sum_{\substack{x_i, y_i \in pairs}} ||Mx_i - y_i||_2^2$$
$$= \arg \min_{M} ||ME^1 - E^2||_2^2$$



- Artetxe, Labake & Agirre, ACL 2017:
 - Use **numbers** as the initial pivot items.
 - Do it an an iterative procedure.

$$\arg \min_{M} \sum_{\substack{x_i, y_i \in pairs}} ||Mx_i - y_i||_2^2$$
$$= \arg \min_{M} ||ME^1 - E^2||_2^2$$



- Artetxe, Labake & Agirre, ACL 2017:
 - Use **numbers** as the initial pivot items.
 - Do it an an iterative procedure.

$$\arg \min_{M} \sum_{\substack{x_i, y_i \in pairs}} ||Mx_i - y_i||_2^2$$
$$= \arg \min_{M} ||ME^1 - E^2||_2^2$$

(also, can be solved exactly with SVD)



Beyond language-to-language

- Words from 1900 to words in 1990
- Words from young speakers to old speakers
- Words from left-wing to right-wing writers
- •

and now for something completely different

Adversarial Examples















speakers













speakers



grasshopper











speakers



grasshopper



poodle



school-bus





...

but also

Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images

Anh Nguyen University of Wyoming anguyen8@uwyo.edu Jason Yosinski Cornell University yosinski@cs.cornell.edu Jeff Clune University of Wyoming jeffclune@uwyo.edu







obelisk





chest









computer keyboard hand blower







assault rifle stethoscope digital clock soccer ball





pinwheel















screwdriver photocopier strawberry

bagel













tile roof ski mask







traffic light





four-poster



African



sea snake





hair slide

nematode



school bus

panpipe



 \boldsymbol{x}

"panda"

57.7% confidence

 $+.007 \times$



sign
$$(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

"nematode"
8.2% confidence



=

$$m{x} + \epsilon \operatorname{sign}(
abla_{m{x}} J(m{ heta}, m{x}, y))$$

"gibbon"
99.3 % confidence

EXPLAINING AND HARNESSING Adversarial Examples

Ian J. Goodfellow, Jonathon Shlens & Christian Szegedy Google Inc., Mountain View, CA {goodfellow, shlens, szegedy}@google.com





Intriguing properties of neural networks

Christian Szegedy	Wojciech Zaremba	Ilya Sutskeve	er Joan Bruna
Google Inc.	New York University	Google Inc.	New York University
Dumitru Erhan	Ian Goodfellow		Rob Fergus
Google Inc.	University of Montreal		New York University
			Facebook Inc.


all of the images on the right are categorized as "ostrich"

Intriguing properties of neural networks

Christian Szegedy Google Inc.	Wojciech Zaremba New York University	Ilya Sutskever Google Inc.	Joan Bruna New York University
Dumitru Erhan	Ian Goodfellow		Rob Fergus
Google Inc.	University of Montreal		New York University

Facebook Inc.

LaVAN: Localized and Visible Adversarial Noise

Danny Karmon¹ Daniel Zoran² Yoav Goldberg¹ ¹ Department of Computer Science, Bar-Ilan University, Ramat Gan, Israel. ² DeepMind, London, UK. karmond@biu.ac.il, danielzoran@google.com, yogo@cs.biu.ac.il

 Original Image Quail: 99.81%, Spiny Lobster: 0.00%
 Noised Image Quail: 0.63%, Spiny Lobster: 94.64%

Quail (99.8%) \rightarrow Spiny Lobster (94.6%)









Drumstick



Milk Can

Printer



Baseball





Fur Coat





Chain Saw





Jaguar



Muzzle



Tiger Cat



Whippet (Dog)







Rock Beauty (Fish)



Tree Frog

Volcano







Vulture



















Padlock (92.7%)

Tiger Cat (94.4%)

Car Mirror (94.5%)

Stingray (90.5%)

Original Image Daisy: 95.64%, Tree Frog: 0.00%





Daisy $(95.6\%) \rightarrow$ **Tree Frog** (96.4%)

Brown Bear: 87.88%, Tree Frog: 0.02%



Noised Image Brown Bear: 5.56%, Tree Frog: 82.69%

Brown Bear (87.9%) \rightarrow Tree Frog (82.7%)





Minivan (90.7%) \rightarrow Tree Frog (86.4%)



how does one find the misleading images?



how does one find the misleading images?

neural network

$$\mathbf{y} = f(\mathbf{x}, \theta)$$

training: gradient decent on parameters to get correct predictions.





training: gradient decent on parameters to get correct predictions.



how does one find the misleading images?

neural network

$$\mathbf{y} = f(\mathbf{x}, \theta)$$

training: gradient decent on parameters to get correct predictions.

adversarial:

change x to get desired y for given params

how does one find the misleading/images?

> neural network $\mathbf{y} = f(\mathbf{x}, \theta)$

training:

gradient decent on parameters to get correct predictions.

adversarial:





how does one find the misleading images?

neural network

$$\mathbf{y} = f(\mathbf{x} + \delta, \theta)$$

training: gradient decent on parameters to get correct predictions.

adversarial:



find small delta to get desired y for given params how does one find the misleading /mages? neural network $\mathbf{y} = f(\mathbf{x} + \delta, \theta)$

training:

gradient decent on parameters to get correct predictions.

adversarial:

Adversarial Examples: Interesting open problems

- How to train a network which is robust to adversarial attacks.
- How to identify adversarial examples?
- How to generate adversarial examples for sequences?

Adversarial Examples: Interesting open problems

- How to train a network which is robust to adversarial attacks.
- How to identify adversarial examples?
- How to generate adversarial examples for sequences?

Adversarial Examples: Interesting open problems

- How to train a network which is robust to adversarial attacks.
- How to identify adversarial examples?
- How to generate adversarial examples for discrete sequences?



Very brief intro: GANs

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

May 18, 2017

Generator network: try to fool the discriminator by generating real-looking images **Discriminator network**: try to distinguish between real and fake images



Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

5

10

Lecture 13 -

Fei-Fei Li & Justin Johnson & Serena Yeung

CS11-747 Neural Networks for NLP Adversarial Methods

Graham Neubig



Carnegie Mellon University

Language Technologies Institute

Site https://phontron.com/class/nn4nlp2017/

Adversarial Feature Learning

Adversaries over Features vs. Over Outputs

Generative adversarial networks



• Adversarial feature learning



- Why adversaries over features?
 - Non-generative tasks
 - Continuous features easier than discrete outputs

Adversaries over Features vs. Over Outputs

Generative adversarial networks



• Adversarial feature learning



Train h to not know something

- Why adversaries over features?
 - Non-generative tasks
 - Continuous features easier than discrete outputs











Adversarial Feature Learning vs Multi-task Learning

(discuss)

Learning Languageinvariant Representations

• Chen et al. (2016) learn language-invariant representations for text classification



Adversarial Language Identification Scorer

Also on multi-lingual machine translation (Xie et al. 2017)



How do we train this in practice?







three different sub-objectives





three different sub-objectives







yellow: update parameters
white: don't update



yellow: update parameters **white**: don't update





yellow: update parameters
white: don't update









Learning Domain-invariant Representations (Ganin et al. 2016)

• Learn features that cannot be distinguished by domain



 Interesting application to synthetically generated or stale data (Kim et al. 2017)
- "Iterative Null-space Projection"
- Goal: a representation h' such that no linear classifier can predict Z from h'.

Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection

Shauli Ravfogel^{1,2} Yanai Elazar^{1,2} Hila Gonen¹ Michael Twiton³ Yoav Goldberg^{1,2} ¹Computer Science Department, Bar Ilan University ²Allen Institute for Artificial Intelligence ³Independent researcher

- "Iterative Null-space Projection"
- Goal: a representation h' such that no linear classifier can predict Z from h'.

more stable than gradient reversal

less heuristic

- "Iterative Null-space Projection"
- Goal: a representation h' such that **no linear** classifier can predict Z from h'.

• Main idea:

train classifier c1. find its null-space. project rep representation to null-space. train classifier c2 on result. find its null-space...



Figure 1: t-SNE projection of GloVe vectors of the most gender-biased words after t=0, 3, 18, and 35 iterations of INLP. Words are colored according to being male-biased or female-biased.



"think about what you create and how."



should we built this?



should we built this?

"online gaming"?

user tracking for oppressive governments?

user tracking for ads?

user tracking for IDF?



should we built this?

"online gaming"?

user tracking for oppressive governments?

user tracking for ads?

user tracking for IDF?

use your own judgement







FACEPTION Facial Personality Analytics



Utilizing advanced machine learning techniques we developed and continue to evolve an array of classifiers. These classifiers represent a certain persona, with a unique personality type, a collection of personality traits or behaviors. Our algorithms can score an individual according to their fit to these classifiers.





score an individual according to their fit to these classifiers.





Who will gain from this? Who will suffer from this?



- Your system will operate in the real world.
- It's decisions will affect real people.
- The fact that something is automated/algorithmic/ learned does not mean it is objective/clean/ harmless/safe
- Someone designed the algorithm.
 Someone created the training data.



algorithms are opinions embedded in code

ML models automate the status quo











some things are obvious, other are more subtle



• The "Gorilla" case.

Google apologises for Photos app's racist blunder

() 1 July 2015







and a state of a

Google apologises for Photos app's racist blunder









Technology

Alex Hern

✓ @alexhern

09.49 BST

old

<

612

Wed 20 May 2015

Chis article is 2 years

93

Flickr faces complaints over 'offensive' auto-tagging for photos

Auto-tagging system slaps 'animal' and 'ape' labels on images of black people, and tags concentration camps with 'jungle gym' and 'sport'



Flickr is facing a user revolt after a new auto-tagging system labelled images of



Why did it happen? How can we treat this?

• [discuss]



Solution? don't say anything is a gorilla

WHEN IT COMES TO GORILLAS, GOOGLE PHOTOS REMAINS BLIND

https://www.wired.com/story/when-it-comes-to-gorillas-google-photos-remains-blind/



ML impacts millions of people everyday

And it also propagates biases in unknown ways





Amazon Reportedly Killed an AI Recruitment System Because It Couldn't Stop the Tool from Discriminating Against Women







by 😔 utbrain 🛛 🕞

27 Gadgets That May Go Out Of Stock Soon by Weekly Penny | Sponsored



Forgetting things at age 40? You should read this! by The Times of Israel | Sponsored







This CIA-funded tool predicts crime before it happens



by BRYAN CLARK — Aug 2, 2017 in ARTIFICIAL INTELLIGENCE

Mos



how will such a system interact with the real world? feedback loops...

T N N



- "Predictive policing"
- Use ML to decide where to send police patrols
 - ...based on previous crime statistics
 - ...reported to the police, and handled by police
- What's the result of this?



- "Predictive policing"
- Use ML to decide where to send police patrols
 - ...based on previous crime statistics
 - ...reported to the police, and handled by police
- What's the result of this?

feedback loops! causality!



- "Predictive policing"
- Use ML to decide where to send police patrols
 - ...based on previous crime statistics
 - ...reported to the police, and handled by police
- What's the result of this?

arrests data != crime data

feedback loops! causality!

- "Crowd-based Video recommendation"
 - YouTube suggests what to watch next based on past user behavior.
 - What are the implications?

- "Crowd-based Video recommendation"
 - YouTube suggests what to watch next based on past user behavior.
 - What are the implications?

- People watch more videos, metrics go up
- but which videos are we watching?
- are we happy with this?



- "Teacher Assessment"
- Teacher is good if many students pass exam.
 - What can happen?

- "Teacher Assessment"
- Teacher is good if many students pass exam.
 - What can happen?
- Teacher is good if many students pass <u>standardized</u> exams.
 - What can happen?



- "Teacher Assessment"
- Teacher is good if many students pass exam.
 - What can happen?
- Teacher is good if many students pass <u>standardized</u> exams.
 - What can happen?
 - factors beyond teacher affect scores
 - influence on teacher's career
 - influence on teacher's choices and the world



Meaning of a mistake

- What can happen if your system is 100% accurate?
- What can happen if your system is 90% accurate?
 - Are the mistakes evenly distributed?



Meaning of a mistake

- A deep learning system that predicts IQ scores based on face images.
 - What if it is 100% accurate?
 - What if it is 90% accurate on females?
 - ... and 60% accurate on blonde females aged < 25?


ML encodes bias

- What biases do you encode into your system?
 - Intentionally
 - Unintentionally
- Who will be affected?



Ethics

- Some possible sources of problems:
 - Bias in training data
 - In p(y)
 - In p(x)

bad data -> bad models

data is based on people...

- In p(y|x)
- In p(x|y)
- Architecture, Features, Loss function
- Bias in outputs
- Feedback loops



Ethics in ML: Some resources

FAT / ML 2018 2017 2016 2015 2014 Organization Resources Mailing list

Fairness, Accountability, and Transparency in Machine Learning

http://www.fatml.org/



Ethics in ML: Some resources

 $\leftarrow \rightarrow$ C (i) Not Secure | demo.clab.cs.cmu.edu/ethical_nlp/

☆ 💟 🕐 🕈 📲 🛋 📄 🚺 🛆

Computational Ethics for NLP

CMU CS 11830, Spring 2018

T/Th 10:30-11:50am, GHC 4215

Yulia Tsvetkov (office hours: Tue noon-1pm, GHC 6405), ytsvetko@cs.cmu.edu Alan W Black (office hours: TBD, GHC 5701), awb@cs.cmu.edu TA: Shrimai Prabhumoye (office hours: Tue 2-3pm, GHC 5511), sprabhum@cs.cmu.edu

Summary Syllabus Readings Grading Projects Policies

http://demo.clab.cs.cmu.edu/ethical_nlp/







Ethics

- Machine learning and Data Science gives great power.
- With great power comes great responsibility.
- Think about what you build and why.

• I hope you enjoyed the course.