### Adversarial Machine Learning Ian Goodfellow, Senior Staff Research Scientist AAAT 2019-01-30

# Google Al



## Most Traditional Machine Learning: Optimization















Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



## Generative Modeling: Sample Generation



Training Data (CelebA)



#### Sample Generator (Karras et al, 2017)



#### (Goodfellow et al., 2014)



### 4.5 years of progress on faces









## 2 Years of Progress on ImageNet





























Odena et al 2016











Miyato et al 2017

Zhang et al 2018

Brock et al 2018

(Goodfellow 2018)

# Unsupervised Image-to-Image Translation





#### Day to night

### (Liu et al., 2017)



# CycleGAN





### (Zhu et al., 2017)



# Video-to-Video

### Pose-to-Body Results











# Everybody Dance Now













## Personalized GANufacturing

#### (Hwang et al 2018)



## Self-Attention





### (Zhang et al., 2018)

Use layers from Wang et al 2018















(Fake)

(Brock et al, 2018)BigGAN Large scale TPU implementation

## Recent Advances

Style-based generators (Karras et al, 2018)



Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



# Adversarial Examples



### $+.007 \times$

#### 58% panda





#### 99% gibbon

## Also Adversarial Examples



(Eykholt et al, 2017)



(Goodfellow 2018)



## Adversarial Examples in the Physical World





(a) Image from dataset

(b) Clean image

(c) Adv. image,  $\epsilon = 4$  (d) Adv. image,  $\epsilon = 8$ 

#### (Kurakin et al, 2016)



## Adversarial Training as a Minimax Problem

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathbb{E}_{\boldsymbol{x},y} \max_{\boldsymbol{\eta}} [J(\boldsymbol{x}, y, \eta)] = (J(\boldsymbol{x}, y, \eta))$$

with the learning algorithm as the minimizing player and a fixed procedure (such as L-BFGS or the fast gradient sign method) as the maximizing player."

- "Adversarial training can be interpreted as a minimax game,
  - $\boldsymbol{\theta}$ ) + J( $\boldsymbol{x} + \boldsymbol{\eta}, \boldsymbol{y}$ )],

- Original implementation: <u>Goodfellow et al 2014</u>
- Explicit use of "minimax": Farley and Goodfellow, 2016





(CleverHans tutorial, using method of Goodfellow et al 2014)





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



# Model-Based Optimization

Make new inventions by finding input that maximizes model's predicted performance







## Designing DNA to optimize protein function





(Gupta and Zou, 2018)

### Make ML work Generative modeling Security Model-based optimization RL

Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



## Adversarial Examples for RL





 $(\underline{\text{Huang et al.}}, 2017)$ 



# Self-Play

#### 1959: Arthur Samuel's checkers agent





#### (OpenAI, 2017)



Goal: push opponent outside the ring, or topple them over

(Bansal et al, 2017)



### SPIRAL Synthesizing Programs for Images Using Reinforced Adversarial Learning

#### Input Program end = [(9, 12), (3, 16), (17, 26), (30, 26), (30, 26), (30, 26), (20, 22), (16, 14), (30, 21), ...], <mark>ctl</mark> = [(8, 11), (8, 24), (3, → 25), (10, 25), (18, 25), (23, 25), (17, 21), (17, 22), (18, 22), ...], pen = [0, 1, 1, 1, 1, 1, 0, Image 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0]

### (Ganin et al, 2018)

#### Interpreters

#### **Simulated Paint**







**Simulated Arm** 

**Real Arm** 





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



- We want extreme reliability for
  - Autonomous vehicles
  - Air traffic control
  - Surgery robots
  - Medical diagnosis, etc.

# Extreme Reliability

• Adversarial machine learning research techniques can help with this

• Katz et al 2017: verification system, applied to air traffic control





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



## Supervised Discriminator for Semi-Supervised Learning





#### (Odena 2016, Salimans et al 2016)

(Goodfellow 2019)

### Virtual Adversarial Training Miyato et al 2015: regularize for robustness to adversarial perturbations of

unlabeled data







Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



### • Domain Adversarial Networks (Ganin et al, 2015)



#### VIPER

• Professor forcing (Lamb et al, 2016): Domain-Adversarial learning in RNN hidden state

# Domain Adaptation

PRID

CUHK



### GANs for simulated training data Unlabeled Real Images







#### Synthetic





#### Refined

(Shrivastava et al., 2016)



# GraspGAN





#### (Bousmalis et al. 2017)





Grasp Success in the Real World

Number of Real-World Samples Used for Training

(Bousmalis et al, 2017)



G

Randomized Simulation





Real World

(James et al, 2018)

## Sim-to-real via sim-to-sim

action

Agent

Agent

action

Canonical Simulation

> Learn to grasp without real data!

Canonical Simulation





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



(Goodfellow 2019)

# Adversarially Learned Fair Representations

- Edwards and Storkey 2015
- Learn representations that are useful for classification
- make S impossible to recover
- Final decision does not depend on S

• An adversary tries to recover a sensitive variable Sfrom the representation. Primary learner tries to



### How do machine learning models work?



(c) Grad-CAM 'Cat'





Interpretability literature: our analysis tools show that deep nets work about how you would expect them to.

(i) Grad-CAM 'Dog' (Selvaraju et al, 2016)



(Goodfellow et al, 2014)

Adversarial ML literature: ML models are very easy to fool and even linear models work in counter-intuitive ways.



### Robust models are more interpretable



Relatively vulnerable model

(Goodfellow 2015)



Relatively robust model





Neuroscience

Fairness, accountability and transparency Domain adaptation

Label efficiency Extreme reliability



### Adversarial examples that affect both computer and time-limited human vision







25% snake

67% snake

Elsayed et al 2018

# Questions

