# Pwning Machine Learning (ML) for Fun and Profit

KIWICON X

Davi Ottenheimer

# Appreciation For Coming

Please Reach Under Your Seat
If You Find A Challenge Token
You Win!

Got Ewe

# Who Defines 'Win'?

- Have You Learned
  - About Chance / Probability
  - About **Authority**
  - About Public Speakers
  - About Me
- 'Control' is a Function of Knowledge

# Security A Competitive Game

1. Out-***Smart*** Adversaries (Preserving Rational Authority)

2. Create Efficiencies, Reduce Costs

3. Balance Both Capabilities!
   - Security **=** ***More*** Success at ***Less*** Cost
   - Don't Confuse/Reverse These

# In a World of
# Decisions For Competition

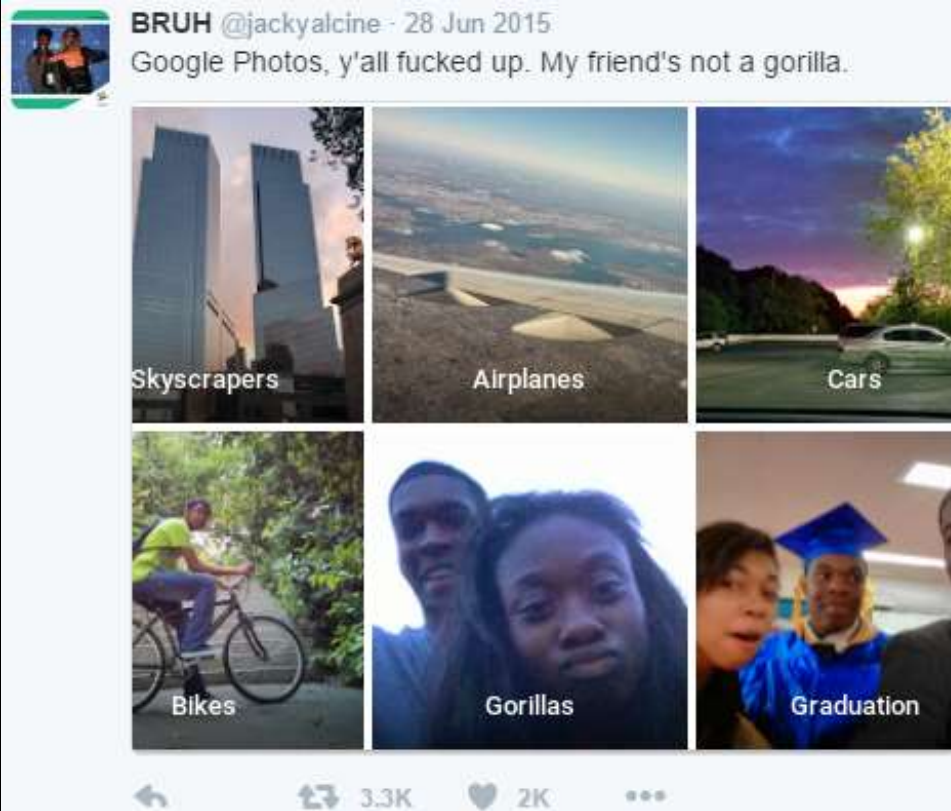We Expect MACHINES
to Make 'Winning' Easier
And Yet...

BAD THINGS ARE HAPPENING

flyingpenguin                    KIWICON X

# 'Really Interesting Problems'

# 'Professional'

# 'Unprofessional'

# False 'Criminal' Labeling

'Blacks falsely labeled future criminals at almost twice the rate of white defendants'

flyingpenguin                          KIWICON X

# False 'Criminal' Labeling

'compared predicted to actual recidivism: scores wrong 40% of the time and biased against black defendants.'

flyingpenguin                                    KIWICON X

# UK False Road Segmentation...

# Death from Autopilot

**60-0 mph Tesla Brake Test = 108 ft**

Chose to kill human because 'overhead sign'
(more likely a *moving* bridge)

https://www.ntsb.gov/investigations/AccidentReports/Pages/HWY16FH018-preliminary.aspx

REASONS FOR
(SAFETY) FAILURES

KIWICON X

# Remember 1958 Predictions?

'The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence'

'You seem a little **shell shocked** so let's just stop you right there...'

flyingpenguin

KIWICON X

# 2016: Bots Are Seeking & Pattern Matching

Algorithms Now Repeat (Faster) Whatever Mistakes They Learn

(see also: decline of civilization)

Image source:

**eseT**
www.eset.com

flyingpenguin

KIWICON X

# Bypass Learning: *Design Flaw*

Neo-Nazis
Piled in a
Backdoor

('Repeat
after me!')

# Break Supervised Learning

- Predicts Future Data Based on Past
- Inputs Data and Labels
  - Classification
  - Regression

| SPAM | NOT SPAM |
|------|----------|

Viagra pills sale          ViAGR4! P1l1s fur sail

# Supervised Break: Face



Detection Result:
Error: "**0 face detected**"

# Supervised Break: Sentiment



Supervised ML
Training Set

Search Algorithm
for Feelings

More Precise Classification
"Very Happy"

| | |
|---|---|
| Very Happy | 0.60 |
| Pain | 0.07 |
| Very Angry | 0.06 |
| Sad | 0.05 |
| Happy | 0.05 |
| Confident | 0.04 |
| Surprised | 0.04 |
| Calm | 0.03 |
| Disgusted | 0.02 |
| Angry | 0.02 |
| Scared | 0.01 |

http://www.datasciencecentral.com/profiles/blogs/tricks-in-face-recognition

flyingpenguin
KIWICON X

# Supervised Break: Fight Faces



Anger 0.01919
Contempt 0.00015
Disgust 0.00192
Fear 0.00817
Happiness 0.91743
Neutral 0.01185
Sadness 0.00780
Surprise 0.03349

"Happiness": 0.917434633,
"Neutral": 0.0118517382,
"Sadness": 0.00780460332,
"Surprise": 0.03348755
}
},
{
"FaceRectangle": {
"Left": 90,
"Top": 32,
"Width": 37,
"Height": 37
},
"Scores": {
"Anger": 0.187864065,
"Contempt": 0.003002729,
"Disgust": 0.0295347776,
"Fear": 0.0174223464,
"Happiness": 0.0223125257,
"Neutral": 0.6334335,
"Sadness": 0.0324874222,
"Surprise": 0.07394262

Face recognition with deep neural networks. http://cmusatyalab.github.io/openface/

# Break Un-Supervised Learning

Discovers Hidden Structures Within Unlabeled Data

- – Compression
- – Clustering



'Ostrich'

# Un-Supervised Break (ID)

Bus        + 'Ostrich'        = 'Ostrich'

90%+ Effective Attack

Christian Szegedy

# Un-Supervised Break (Traffic)

No Parking          + 'Stop'            = 'Stop'

# Un-Supervised Break (Lip Read)

Confused by English Accent



(a) Lip-rounding vowels    (b) Alveolar    (c) Bilabial    (d) Viseme Categories

Figure 3: Intra-viseme and inter-viseme confusion matrices, depicting the three categories with the most confusions, as well as the confusions between viseme clusters. Colours are row-normalised to emphasise the errors.

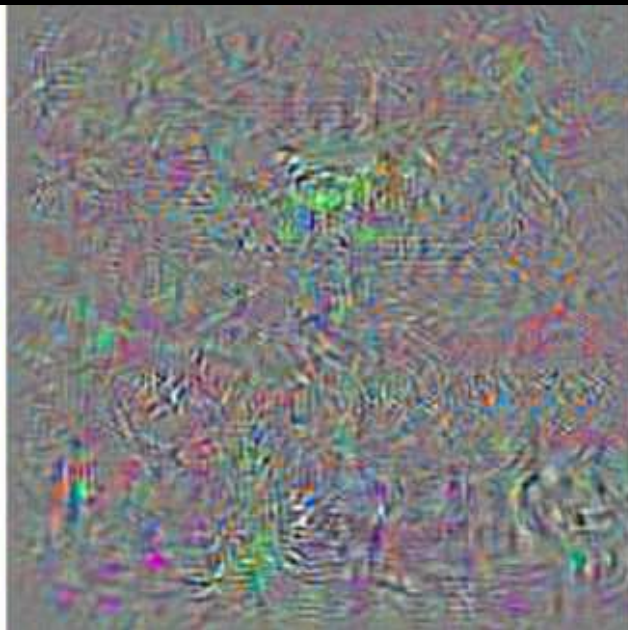Given that the speakers are British, the confusion between /aa/ and /ay/ (Figure 3a) is most probably due to the fact that the first element, and the greater part, of the diphthong /ay/ is articulatorily identical with /aa/: an open back unrounded vowel (Ferragne & Pellegrino, 2010). The confusion of /ih/ (a rather close vowel) and /ae/ (a very open vowel) is at first glance surprising, but in fact in the sample /ae/ occurs only in the word *at*, which is a function word normally pronounced with a reduced, weak vowel /ah/. /ah/ and /ih/ are the most frequent unstressed vowels and there is a good deal of variation within and between them, e.g. *private* and *watches* (Cruttenden, 2014).

LipNet, Sentence-Level Lipreading, Nov 2016
https://arxiv.org/pdf/1611.01599v1.pdf

flyingpenguin                    KIWICON X

# Break Reinforcement Learning

- Improves Performance in Dynamic Environs Using Delayed Rewards

- Measures Actions to Achieve Set Goals (Win a Game)



KASPAROV BEATS 'DEEP BLUE' IN ONE MOVE

OFF

# Reinforcement Learning Defeat



What Do You See?

Lepse Avenue, Kiev, Ukraine: https://www.jwt.com/en/ukraine/work/pedestrianghost/

# Reinforcement Learning Defeat

# Reinforcement Learning Defeat

TO FIX ML SAFETY,
WE NEED TO *BREAK ML MORE*

flyingpenguin                    KIWICON X

Cognitive Bias Options

(Menu of Attacks)

LIMIT TO MEMORY

TOO MUCH INFO

NEED TO ACT FAST

NOT ENOUGH MEANING

To avoid mistakes, we aim to preserve autonomy and group status, and avoid irreversible decisions

To get things done, we tend to complete things we've invested time & energy in

To stay focused, we favor the immediate, relatable thing in front of us

To act, we must be confident we can make an impact and feel what we do is important

We tend to find stories and patterns even when looking at sparse data

We fill in characteristics from stereotypes, generalities, and prior histories

We imagine things and people we're familiar with or fond of as better

We simplify probabilities and to make them easier to think

We think we know what other people are thinking

We project our current mindset and assumptions onto the past and future

flyingpenguin     KIWICON X

# Machines Sort of Like Humans

Expecting machine intelligence
to evolve a human-like brain
like
waiting for airplanes
to grow feathers

# Human Learning as Philosophy

Rene Descartes
(1596-1650)



**1637:** 'Cogito, ergo sum'

John Locke
(1632-1704)



**1693:** Reflective Process,
Articulated Steps

# False Win Can Mean 'Not False' Despite *All Models Being Wrong*

A 'winning result' is not yet proven enough to be flagged wrong or ruled out*

BRUH @jackyalcine · 28 Jun 2015
Google Photos, y'all fucked up. My friend's not a gorilla.

* Transparency of ML may be inversely related to accountability or legal liability

# Machine 'Reflective Process'
## Expressed as Holdout Method

| Initial Data | | |
|---|---|---|

| Evaluate Algorithms | Testing | Training | |
|---|---|---|---|

| Re-Evaluate Algorithms | Testing | Training | Validation | Model Tune & Evaluate |
|---|---|---|---|---|

ML Algorithm

Model

# Bias is Unsafe:
## Powerful Man Promises A Great Opportunity!

# Bias is Unsafe:
## This is Fine (TIF)

# Opportunity + TIF = Tesla
## (Cognitive Errors)

- AutoPilot Thought It Was Winning
  - Need to Act Fast
  - Make An Impact
- Human Thought He Was Winning
  - Imagined Tesla Better
  - Not Enough Meaning
- ML Opacity + False Wins Led to Death

BUT WAIT...IT GETS WORSE

# The Inevitable Militarization of [Insert Technology Here]*

*Artificial Intelligence
http://www.cyberdefensereview.org/2016/02/08/the-inevitable-militarization-of-artificial-intelligence/

# History



Maxim, an egomaniacal draft dodger, gave the world the first true automatic weapon (Patent No 3493 1883). Used by British in Colonial Africa and by Germans in WWI to *turn earth into hell*. Died proud.

– C. J. Chivers

# Harvard Math PhD is Worried

# The O'Neil Guide to
## Algorithms We Should Worry About

| Widespread Impact | Secret (Targets Can't Understand) | Destructive (Ruin Life, Unfair) |
|---|---|---|

# The O'Neil Guide to
## Algorithms We Should **NOT** Worry About

flyingpenguin

KIWICON X

# NetFlix is Watching You

- Widespread
- Secret
- Destructive?
  - Content Type
  - View Time
  - Correlation

# Widespread and Secret Cameras
## Story-Driven Summarization for Egocentric Video



Figure 6. Example from UTE data comparing our summary (top row) to the three baselines. Our method clearly captures the progress of the story: serving ice cream leads to weighing the ice cream, which leads to watching TV in the ice cream shop, then driving home. Even when there are no obvious visual links for the story, our method captures visually distinct scenes (see last few subshots in top row). The shortest-path approach makes abrupt hops across the storyline in order to preserve subshots that smoothly transition (see redundancy in its last 5 subshots). While the object-driven method [14] does indeed find some important objects (e.g., TV, person), the summary fails to suggest the links between them. Note that object-driven method sometimes produces shorter summaries (like this example) depending on number of unique important objects discovered in the video. See supplementary file for videos.

**US Drones in 2009 generated 24 Years of Video**

http://www.cs.utexas.edu/~grauman/papers/lu-grauman-cvpr2013.pdf, The Economist "The Data Deluge"

# Widespread and Secret Health Sensors

'we know the estimated numbers of people being served by each waste water treatment plant, we can back-calculate daily loads'

Analysis of Greater Chicago Wastewater **(1.5 billion gal/day)**

- Disease
- Drugs
- Environmental Risk

http://gizmodo.com/5844925/chicagos-stickney-wastewater-treatment-plant-is-the-crappiest-place-on-earth,
http://planetearth.nerc.ac.uk/news/story.aspx?id=1185
http://www.treehugger.com/natural-sciences/fish-near-water-treatment-plants-are-harmed-by-human-drugs.html

# Widespread and Secret Drug Sensors

Croatia
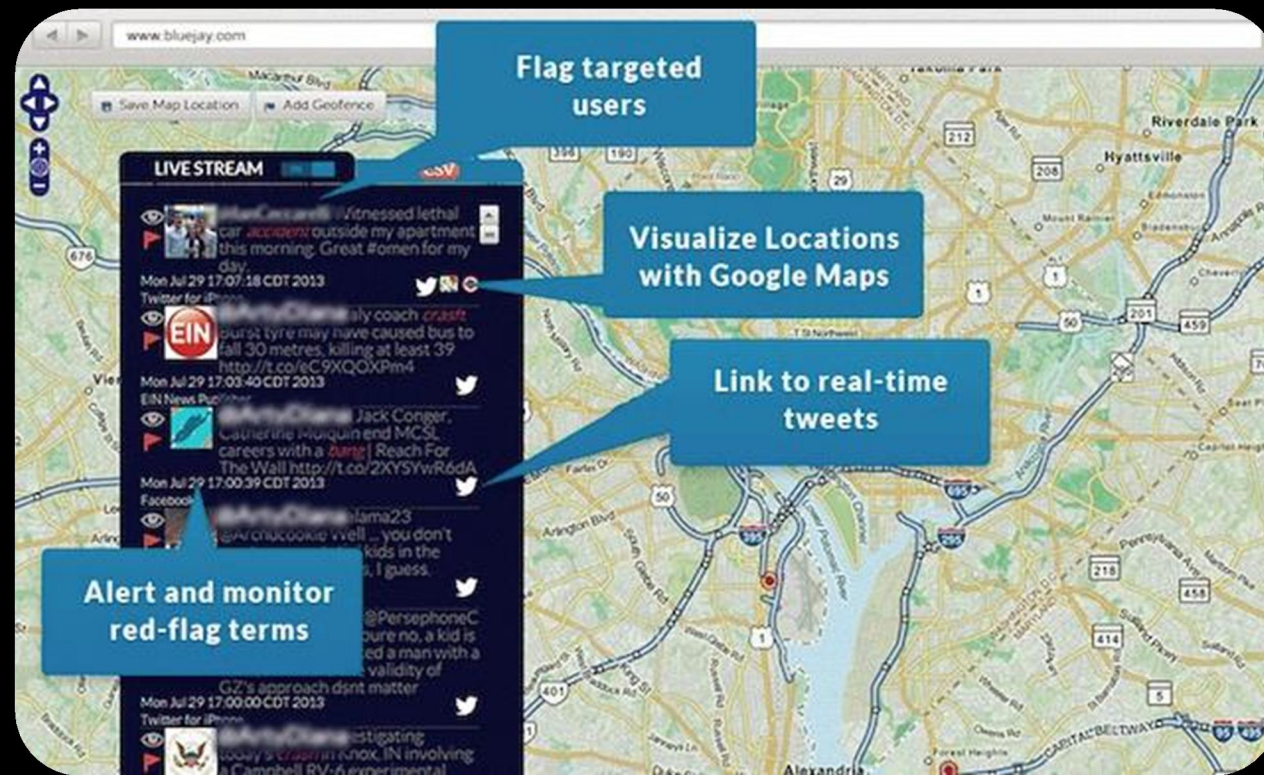
Italy

Finland

London

Oregon

Canada

'Cocaine so widely used it has contaminated Britain's drinking water'

http://gizmodo.com/meth-in-london-heroin-in-zagreb-the-answer-is-found-i-1508209127
http://www.thesundaytimes.co.uk/sto/news/uk_news/Health/article1409450.ece

flyingpenguin                                KIWICON X

# Widespread and Secret Police Sensors
## From BlueToad to BlueJay

# Facebook's Hunter Killer Project
## 'Trajectory Targeting by Prediction'

The goal of this competition is to predict which place a person would like to check in to. For the purposes of this competition, Facebook created an artificial world consisting of more than 100,000 places located in a 10 km by 10 km square. For a given set of coordinates, your task is to return a ranked list of the most likely places. Data was fabricated to resemble location signals coming from mobile devices, giving you a flavor of what it takes to work with real data complicated by inaccurate and noisy values. Inconsistent and erroneous location data can disrupt experience for services like Facebook Check In.

https://www.kaggle.com/c/facebook-v-predicting-check-ins

flyingpenguin                                    KIWICON X

# PWN ML FOR FUN & PROFIT
## (SUPPORTING MUTUAL INTERESTS)*
# OR GET PWNED...

* Thomas Hobbes

flyingpenguin

Pwning Machine Learning (ML)
for Fun and Profit

KIWICON X

Davi Ottenheimer

flyingpenguin

# EPILOGUE

# Why Nazis Love Saying 'Fake'

- Refusal to learn is a 'deliberate, often psychologically motivated, neglect of information too upsetting to allow'

- 'Never believe that anti-Semites are completely unaware of the absurdity of their replies. They know that their remarks are frivolous, open to challenge.' - Jean Paul Sartre

http://www.washingtonpost.com/sf/brand-connect/bleecker-street/denial/
http://abahlali.org/files/Jean-Paul_Sartre_Anti-Semite_and_Jew_An_Exploration_of_the_Etiology_of_Hate__1995.pdf

# Why Nazis Love Saying 'Fake'

'...they are amusing themselves, for it is their adversary who is obliged to use words responsibly, since he believes in words. The anti-Semites have the right to play. They even like to play with discourse for, by giving ridiculous reasons, they discredit the seriousness of their interlocutors. They delight in acting in bad faith, since they seek not to persuade by sound argument but to intimidate and disconcert. If you press them too closely, they will abruptly fall silent, loftily indicating by some phrase that the time for argument is past. It is not that they are afraid of being convinced. They fear only to appear ridiculous or to prejudice by their embarrassment their hope of winning over some third person to their side.'

– Jean Paul Sartre

flyingpenguin                    KIWICON X