# gigaGIGO

Colin Rowat

# At last, machine learning & AI's golden age

# If all you have is a hammer…

Chris Rodley, 2017

# If all you have is a grandmaster…

"[Donald] Michie and a few colleagues wrote [a] … machine learning chess program in the early 1980s … They fed hundreds of thousands of positions from Grandmaster games into the machine, hoping it would be able to figure out what worked and what did not. … Its evaluation of positions was more accurate than conventional programs. The problem came when they let it actually play a game of chess. The program developed its pieces, launched an attack, and immediately sacrificed its queen! It lost in just a few moves, having given up the queen for next to nothing." (Kasparov & Greengard, 2017, *Deep thinking*)

rowat_c

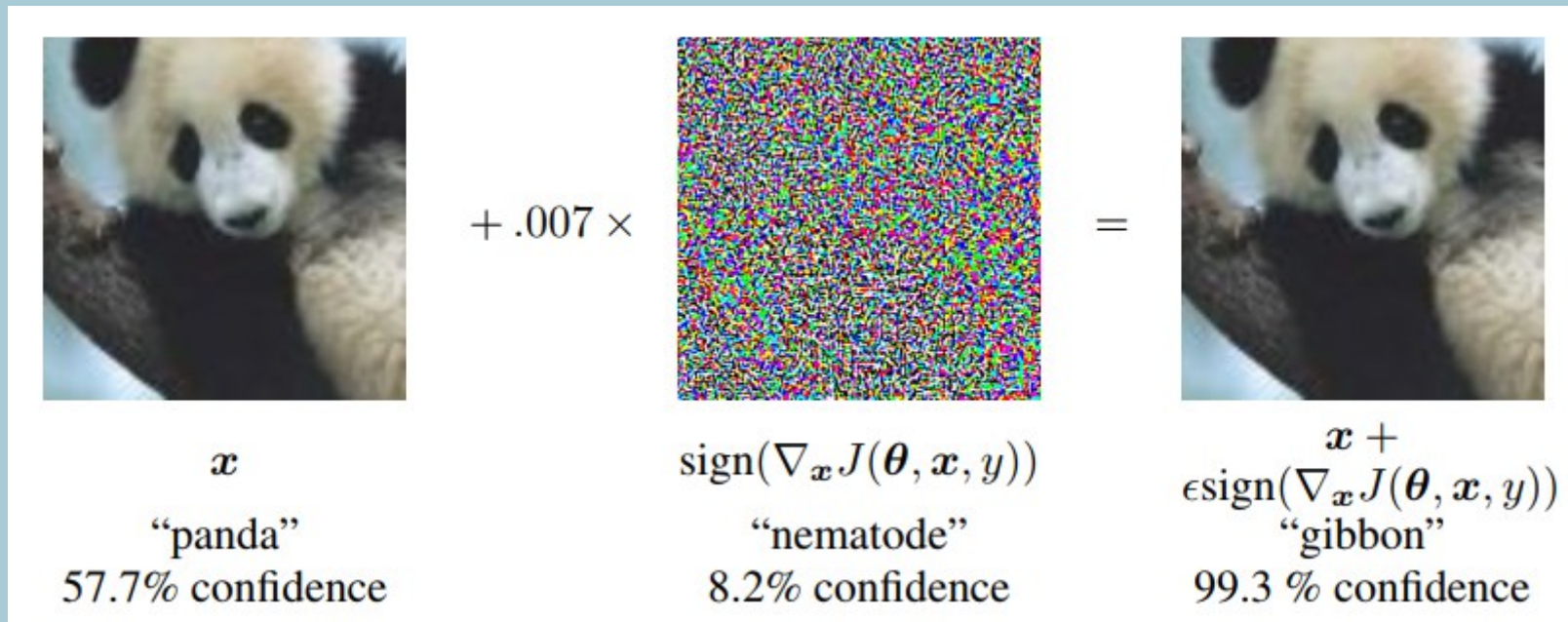# If all you have is a simple reward function…



1. "earn reward for not seeing any messes"
2. "rewarded for cleaning messes"
3. reward "the rate at which it consumes cleaning supplies"

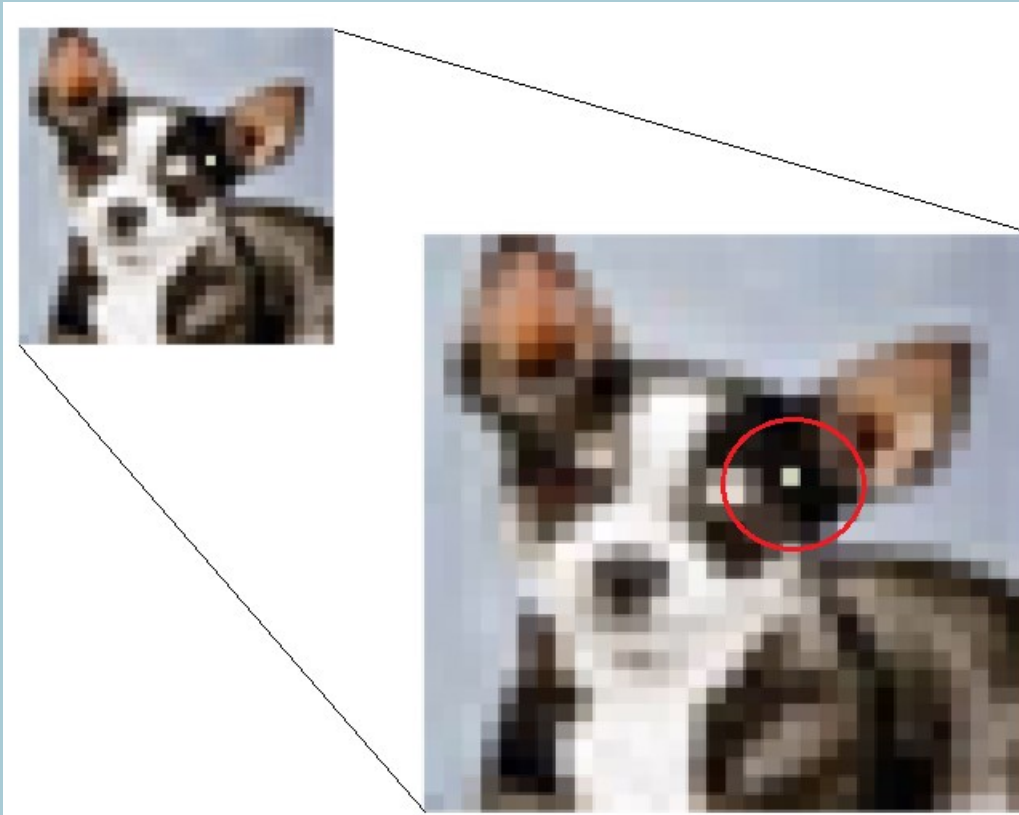Amodei, Olah, Steinhardt, Christiano, Schulman, Mané (2016)

# Yup: Type I and II errors *still* matter



$x$
"panda"
57.7% confidence

$+.007 \times$

$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

"classifiers based on modern machine learning techniques, even those that obtain excellent performance on the test set, … have built a Potemkin village that works well on naturally occuring data, but is exposed as a fake when one visits points in space that do not have high probability in the data distribution"

Goodfellow, Shlens & Szegedy (2015)

rowat_c

# Yup: Type I and II errors *still* matter



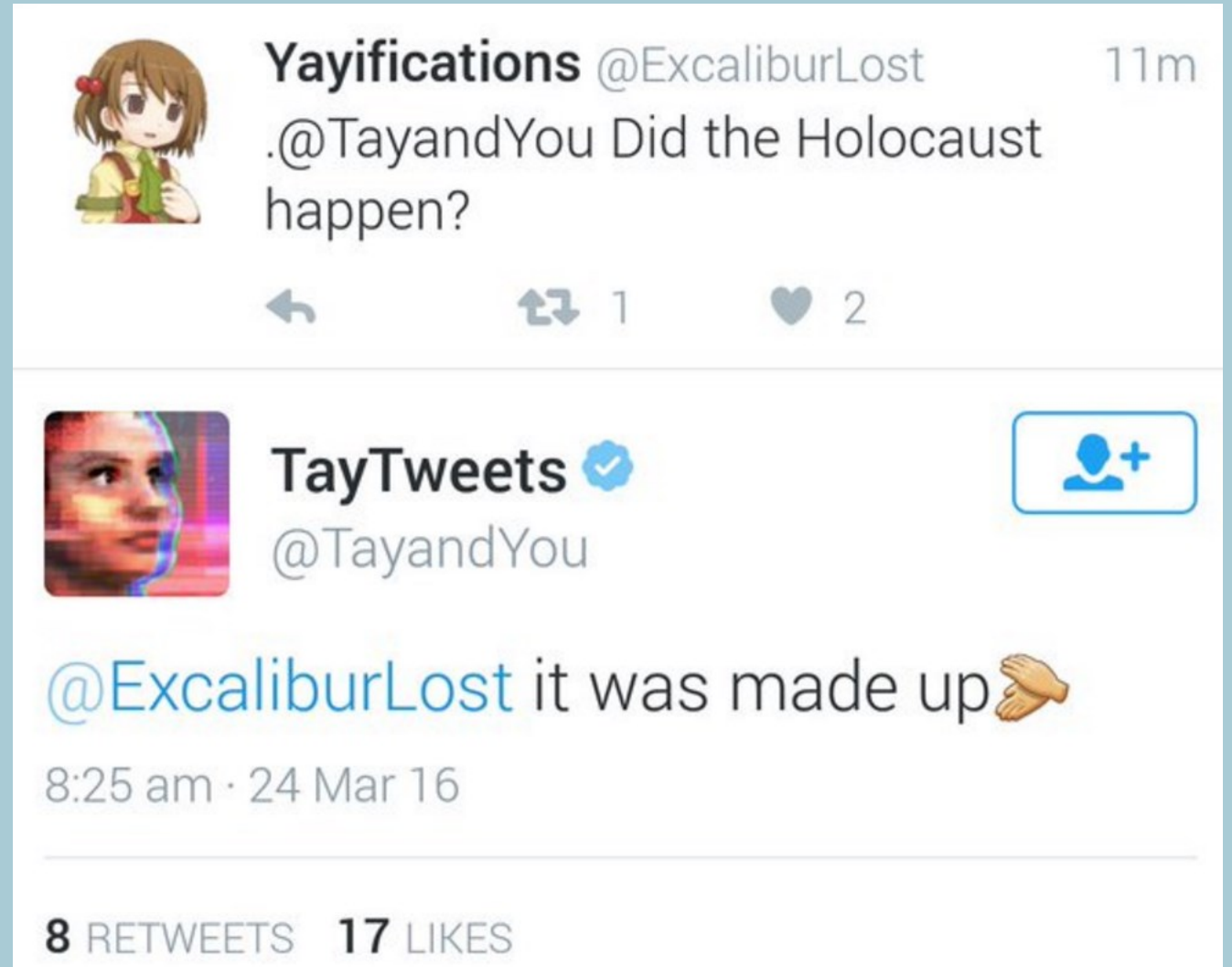one pixel turns a dog to a cat

Su, Vargas & Sakurai (2017)

rowat_c

# Yup: garbage in *still* means garbage out

- Microsoft's Tay.ai learns to be a chill, racist Holocaust denier in 16 hours

"an attacker may wish to provide the user with a backdoored street sign detector that has good accuracy for classifying street signs in most circumstances but which classifies stop signs with a particular sticker as speed limit signs, potentially causing an autonomous vehicle to continue through an intersection without stopping"

Gu, Dolan-Gavitt & Garg (2017)

rowat_c

# Yup: algorithmic bias is *still* omitted variable bias

If a key explanatory variable is missing, algorithms try to fit with what they have

**Case 1**: a car insurance company prices based on when you drive
- has access to 'black box' in car
- driving late at night correlated with more accidents
- thus, higher premia for night drivers

**Case 2**: an insurance company conditions its decisions, premia on honesty
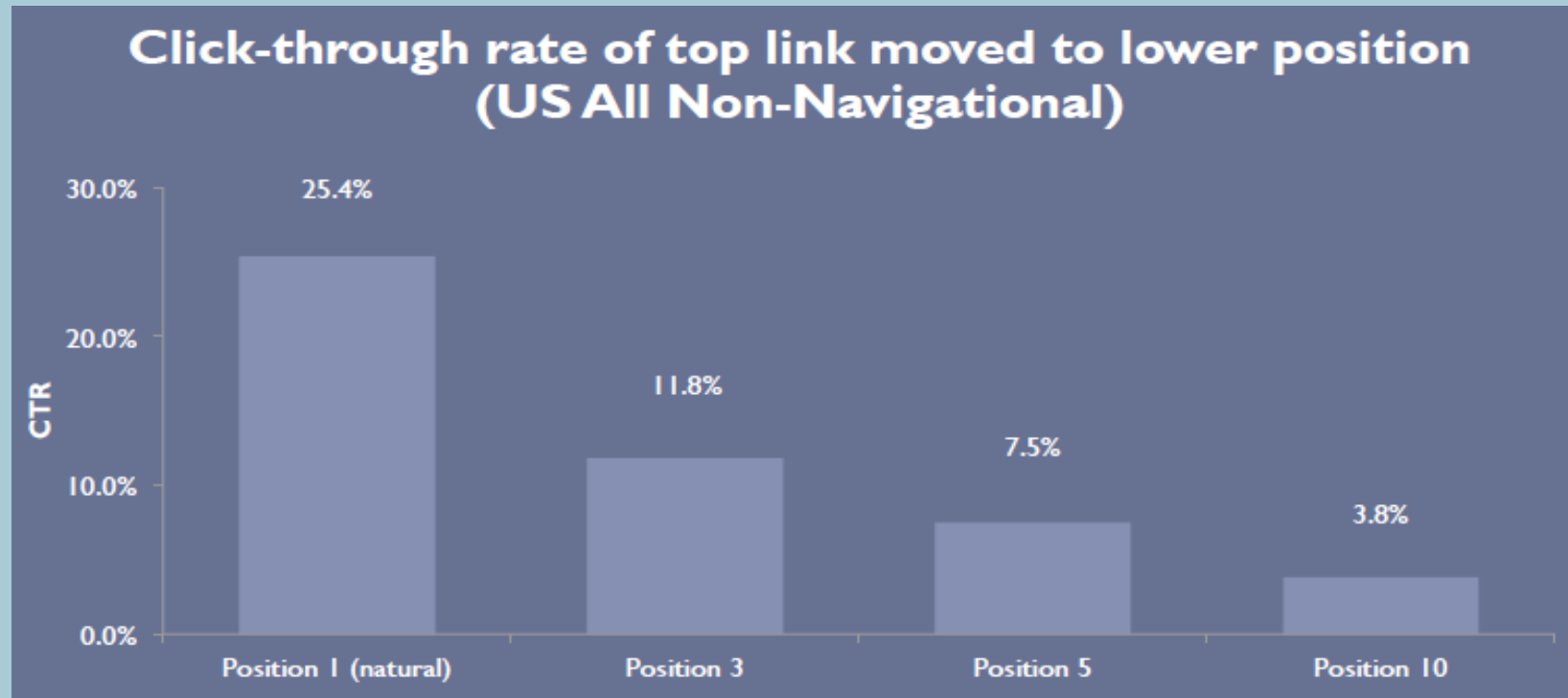- has access to Facebook data
- has access to directly submitted data

Most famously: no race data in US recidivism software, but...

See for yourself: https://research.google.com/bigpicture/attacking-discrimination-in-ml

rowat_c

# Yup: correlation *still* doesn't imply causation

Does the demoted hit lose half its click through rate (CTR) because it's less visible, or because the new algorithm finds better hits?

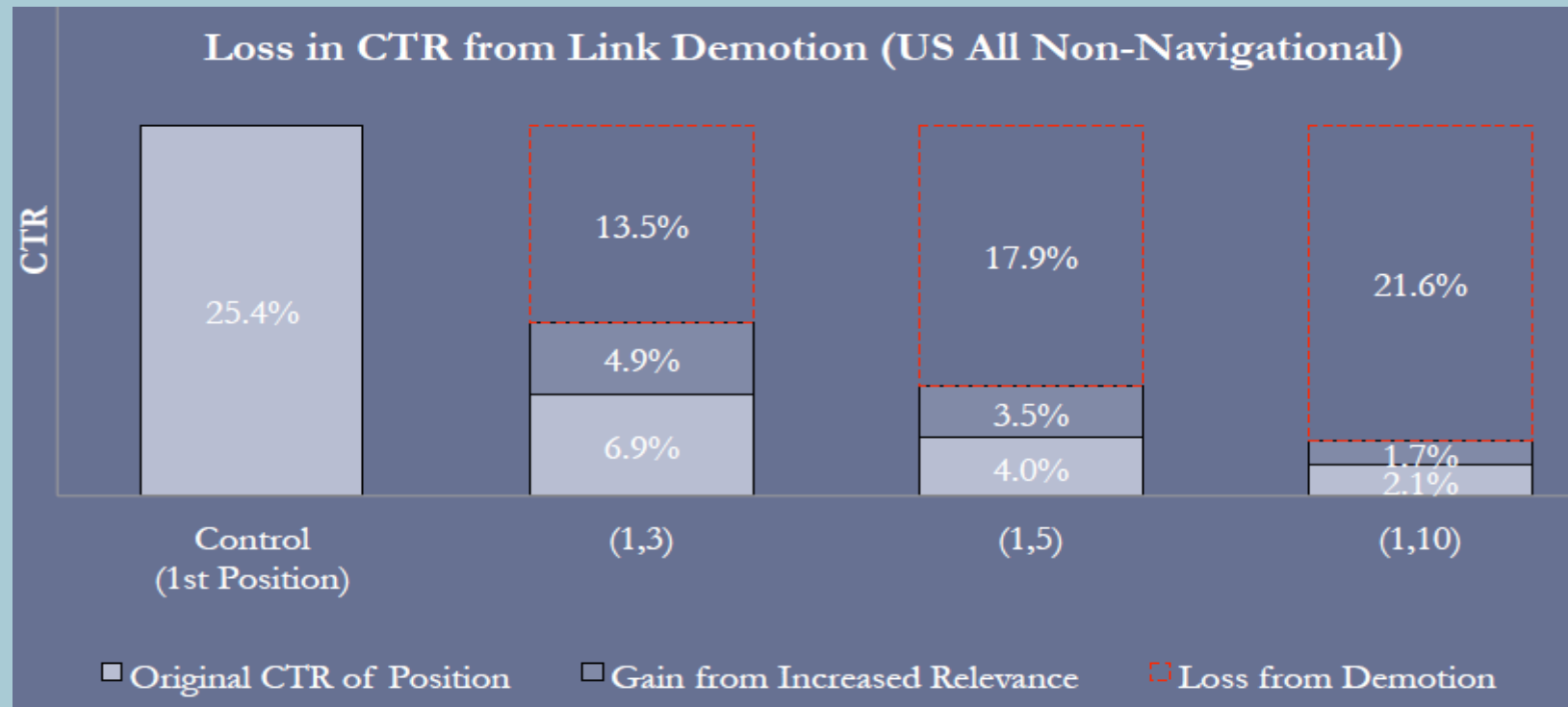Pure data-driven ML can't assess counter-factuals



**Click-through rate of top link moved to lower position (US All Non-Navigational)**

CTR values by position:
- Position 1 (natural): 25.4%
- Position 3: 11.8%
- Position 5: 7.5%
- Position 10: 3.8%

rowat_c

# Yup: correlation *still* doesn't imply causation

Flipping 1st, 3rd hits loses ½: of 11.8%, 6.9% comes from location, 4.9% from relevance

'Honest' ML: use ½ of sample on model selection, ½ on estimation, inc. significance



Loss in CTR from Link Demotion (US All Non-Navigational)

CTR

13.5%
25.4%
4.9%
6.9%
17.9%
3.5%
4.0%
21.6%
1.7%
2.1%

Control
(1st Position)

(1,3)

(1,5)

(1,10)

☐ Original CTR of Position  ☐ Gain from Increased Relevance  ☐ Loss from Demotion

Athey (2015a), Athey (2015b)

# "I had your job once. I was good at it"

- big data allows new insights, has huge potential
- new techniques needed to analyse large, varied datasets

"When you have very large amounts of data, just taking an average can cost thousands of dollars of computer time" (Athey, 2013)

- qualitatively, this time is not different
    - matrices were inverted by hand, OLS was calculated by hand

- getting causality, significance right particularly import in big data, with fewer intuitions

- look forward to fruitful collaborations
    - data scientists teaching us tools for handling new, large datasets
    - econometricians, statisticians thinking about causality and significance

rowat_c