# Common Pitfalls for Studying the Human Side of Machine Learning

**Joshua A. Kroll**, **Nitin Kohli**, Deirdre Mulligan
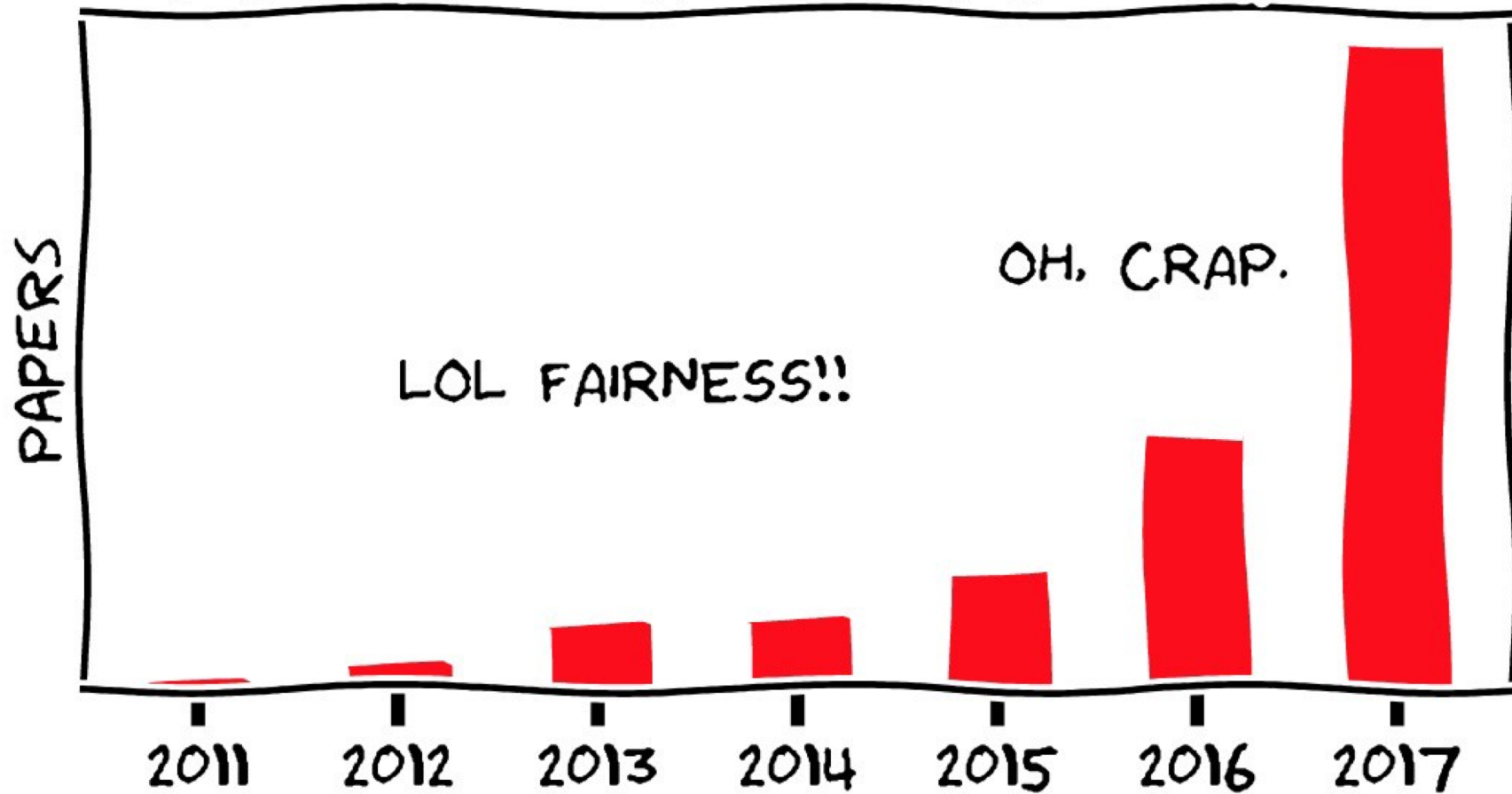
UC Berkeley School of Information
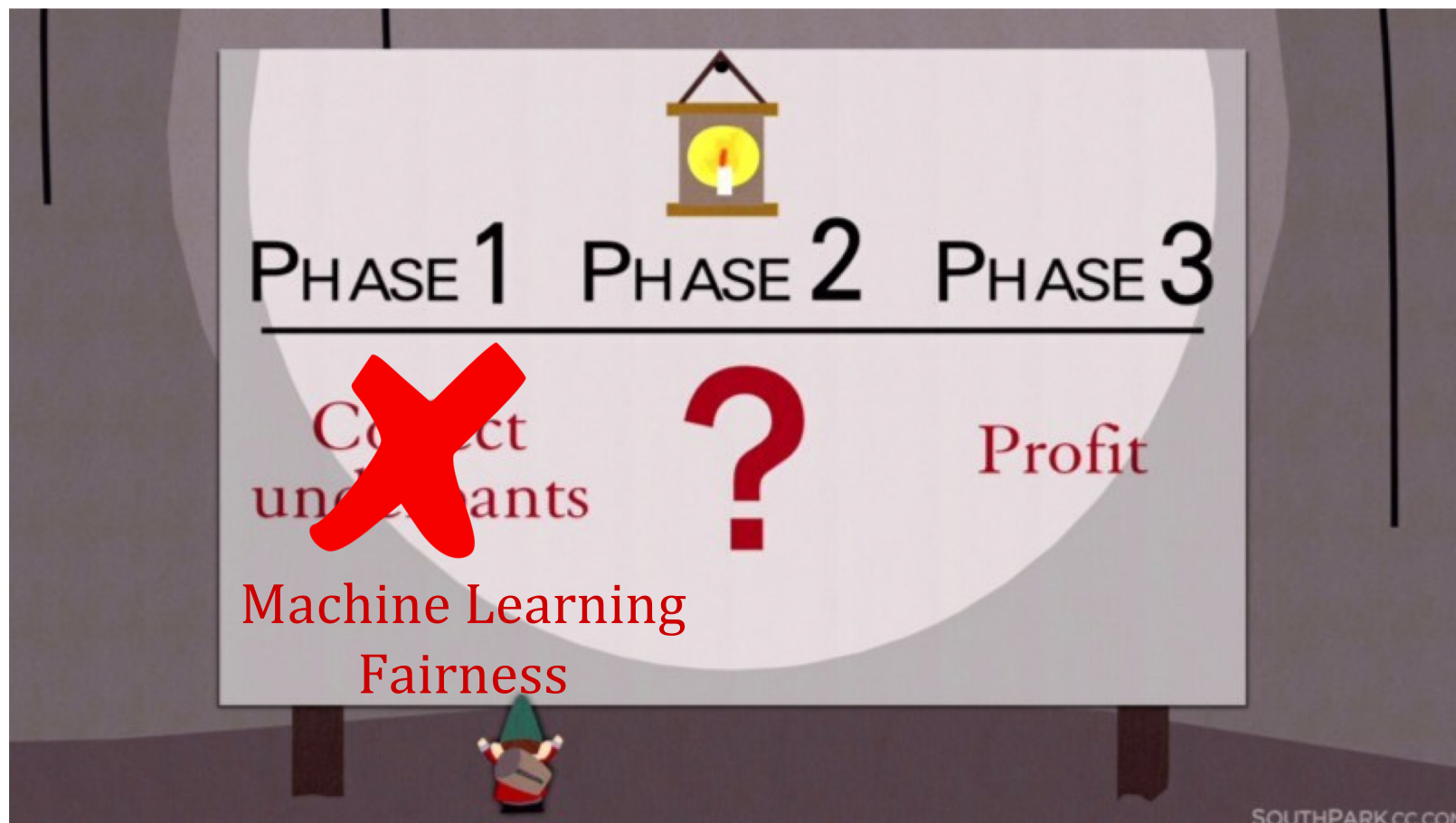
Tutorial: NeurIPS 2018

3 December 2018

Berkeley SCHOOL OF INFORMATION

Credit: Last Year, Solon Barocas and Moritz Hardt, "Fairness in Machine Learning", NeurIPS 2017

# What goes wrong when engaging other disciplines?

- Want to build technology people can *trust* and which supports *human values*
- Demand for:
  - **Fairness**
  - **Accountability**
  - **Transparency**
  - **Interpretability**
- These are rich concepts, with long histories, studied in many ways


- But these terms get re-used to mean different things!

  - This causes unnecessary misunderstanding and argument.

  - We'll examine different ideas referenced by the same words, and examine some concrete cases

# Why this isn't ethics

Machine learning is a tool that solves specific problems

Many concerns about computer systems arise not from people being unethical, but rather from misusing machine learning in a way that clouds the problem at hand

Discussions of ethics put the focus on the individual actors, sidestepping social, political, and organizational dynamics and incentives

Definitions are unhelpful
(but you still need them)

# Values Resist Definition

Definitions aren't for everyone:
Where you sit is where you stand

If we're trying to capture human values,
perhaps mathematical correctness isn't enough

These problems are *sociotechnical* problems

# Fairness

*"What is the problem to which fair machine learning is the solution?" - Solon Barocas*

# What is Fairness:
# Rules are not processes

Tradeoffs are inevitable

# Maybe the Problem is Elsewhere

# What is Accountability: Understanding the Unit of Analysis

What should be true of a system, and where should we intervene on that system to guarantee this?

| | | |
|---|---|---|
| | | 500 |
| | | 500 |
| | | 40 |
| | | 363 |
| | | 100 |
| | | 18 |
| | | 13 |
| | | 25 |
| | | 1.34 |
| | | 12.35 |
| | | 4.30 |
| | | 4.2 |
| | | 16 |
| | | 4.31 |
| | | 304 |
| | | 1.98 |
| | | 11 |
| | | 16 |
| | | 1.20 |
| | | 2.8 |
| | | 1050.1 |

Right page:

| | |
|---|---|
| | 25 |
| | 2 |
| | 27 |
| Summe | 1050.1 |
| | 27.3 |
| | 1022.33 |

# Transparency & Explainability are Incomplete Solutions

# Transparency

| | Cleaned repo | Latest commit `cfc2205` on Sep 21 |
|---|---|---|
| .. | | |
| 📄 baseline.py | Cleaned repo | 2 months ago |
| 📄 cnn-feats-svm.py | Cleaned repo | 2 months ago |
| 📄 cnn.py | Cleaned repo | 2 months ago |
| 📄 decisiontree.py | Cleaned repo | 2 months ago |
| 📄 extract-cnn-feats.py | Cleaned repo | 2 months ago |
| 📄 logistic.py | Cleaned repo | 2 months ago |
| 📄 lstm.py | Cleaned repo | 2 months ago |
| 📄 majority-voting.py | Cleaned repo | 2 months ago |
| 📄 maxent-nltk.py | Cleaned repo | 2 months ago |
| 📄 naivebayes.py | Cleaned repo | 2 months ago |
| 📄 neuralnet.py | Cleaned repo | 2 months ago |

# Explainability

# Explanations from Miller (2017)

- Causal
- Contrastive
- Selective
- Social
- Both a product and a process

Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences."
arXiv preprint arXiv:1706.07269 (2017).

# Data are not the truth

FRANCE

PORTUGAL

MADRID

Barcelona

Palma

Murcia

Sevilla

Gibraltar

Ceuta

Melilla

ALGERIA

MOROCCO

www.comersis.com

If length is hard to measure,
what about unobservable
constructs like risk?

# Construct Validity

Abstraction is a fiction

There is no substitute for
solving the problem

You must first understand
the problem

# Case One :
# Babysitter Risk Rating

Xcorp launches a new service that uses social media data to predict whether a babysitter candidate is likely to abuse drugs or exhibit other undesirable tendencies (e.g. aggressiveness, disrespectfulness, etc.)

Using computational techniques, Xcorp will produce a score to rate the riskiness of the candidates. Candidates must opt in to being scored when asked by a potential employer.

This product produces a rating of the quality of the babysitter candidate from 1-5 and displays this to the hiring parent.

With a partner, examine the validity of this approach. Why might this tool concern people, and who might be concerned by it?

# What would it mean for this system to be fair?

# What would we need to make this system sufficiently transparent?

Are concerns with this system solved by explaining outputs?

# Possible solutions?

This is not hypothetical.

**The Washington Post**
*Democracy Dies in Darkness*

The Switch

# Wanted: The 'perfect babysitter.' Must pass AI scan for respect and attitude.

By **Drew Harwell**
November 23

Read more here:

https://www.washingtonpost.com/technology/2018/11/16/wanted-perfect-babysitter-must-pass-ai-scan-respect-attitude/

(Break)

# Case Two:
# Law Enforcement Face Recognition

The police department in Yville wants to be able to identify criminal suspects in crime scene video to know if the suspect is known to detectives or has been arrested before.

Zcorp offers a cloud face recognition API, and the police build a system using this API which queries probe frames from crime scene video against the Yville Police mugshot database.

What does the fact that this is a government application change about the requirements?

What fairness equities are at stake in such a system?

# What is the role of transparency here?

Who has responsibility in or for this system?
What about for errors/mistakes?

What form would explanations take in this system?

This is not hypothetical, either.

# Amazon's Face Recognition Falsely Matched 28 Members of Congress With Mugshots

By Jacob Snow, Technology & Civil Liberties Attorney, ACLU of Northern California
JULY 26, 2018 | 8:00 AM

TAGS: Face Recognition Technology, Surveillance Technologies, Privacy & Technology

Read more here:
https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28

To solve problems with machine learning, you must understand them

Respect that others may
define the problem
differently

If we allow that our systems include people and society, it's clear that we have to help **negotiate** values, not simply **define** them.

There is no substitute for thinking

Questions?