



The Impossibility of Productizable AI: Problems and Potential Solutions

Falaah Arif Khan,
Identity and Payments Engg, Dell Commerce Platform

Personal: <https://falaaharifkhan.github.io/research>

Blog: <https://thefaladox.wordpress.com/>

Twitter: @ArifFalaah



To Do:

- ☐ What are the problems?
- ☐ How do these problems arise?
- ☐ Why should we care?
- ☐ How do we fix these problems?





FACEBOOK'S ARTIFICIAL INTELLIGENCE ROBOTS SHUT DOWN AFTER THEY START TALKING TO EACH OTHER IN THEIR OWN LANGUAGE

'you i i i everything else'

Andrew Griffin | @_andrew_griffin | Monday 31 July 2017 17:10 |
88 comments



07-14-17

AI Is Inventing Languages Humans Can't Understand. Should We Stop It?

Researchers at Facebook realized their bots were chattering in a new language. Then they stopped it.

“Agents will drift off understandable language and invent codewords for themselves,” says Batra, speaking to a now-predictable phenomenon that’s been observed again, and again, and again. “Like if I say ‘the’ five times, you interpret that to mean I want five copies of this item. This isn’t so different from the way communities of humans create shorthands.”



Facebook put cork in chatbots that created a secret language

Alice and Bob, the two bots, raise questions about the future of artificial intelligence.



Richard Nieva July 31, 2017 11:58 AM PDT

ES



A pair of chatbots has recently done something children often do: create a secret language.

Last month, researchers at Facebook found two bots developed in the social network's AI division had been communicating with each other in an unexpected way. The bots, named Bob and Alice, had generated a language all



Divide these objects between you and another Turker. Try hard to get as many points as you can!

Send a message now, or enter the agreed deal!

Items	Value	Number You Get
	8	<input type="text" value="1"/>
	1	<input type="text" value="1"/>
	0	<input type="text" value="0"/>

Fellow Turker: I'd like all the balls

You: Ok, if I get everything else

Fellow Turker: If I get the book then you have a deal

You: No way - you can have one hat and all the balls

Fellow Turker: Ok deal

Type Message Here:

Figure 1: A dialogue in our Mechanical Turk interface, which we used to collect a negotiation dataset.

Reinforcement learning for dialog agents

Supervised learning aims to imitate the actions of human users, but it does not explicitly attempt to achieve an agent's goals. Taking a different approach, the FAIR team explored pre-training with supervised learning, and then fine-tuned the model against the evaluation metric using reinforcement learning. In effect, they used supervised learning to learn how to map between language and meaning, but used reinforcement learning to help determine which utterance to say.

During reinforcement learning, the agent attempts to improve its parameters from conversations with another agent. While the other agent could be a human, FAIR used a fixed supervised model that was trained to imitate humans. The second model is fixed, because the researchers found that updating the parameters of both agents led to divergence from human language as the agents developed their own language for negotiating. At the end of every dialog, the agent is given a reward based on the deal it agreed on. This reward was then back-propagated through every word that the agent output, using policy gradients, to increase the probability of actions that lead to high rewards.

Models produce meaningful novel sentences.

One interesting question is whether our models are capable of generating novel sentences in the new circumstances they find themselves in, or if they simply repeat messages from the training data verbatim. We find that 76% of messages produced by the LIKELIHOOD model in self-play were found in the training data. We manually examined the novel utterances produced by our model, and found that the overwhelming majority were fluent English sentences in isolation—showing that the model has learnt a good language model for the domain (in addition to results that show it uses language effectively to achieve its goals). These results suggest that although neural models are prone to the safer option of repeating sentences from training data, they are capable of generalising when necessary. Future work should choose domains that force a higher degree of diversity in utterances.

Maintaining multi-sentence coherence is challenging. One common linguistic error we see RL+ROLLOUTS make is to start a message by indicating agreement (e.g. *I agree* or *Deal*), but then going on to propose a counter offer—a behaviour that human partners found frustrating. One explanation is that the model has learnt that in the supervised data, messages beginning with *I agree* are often at the end of the dialogue, and partners rarely reply with further negotiation—so the models using rollouts and reinforcement learning believe this tactic will help their offer to be accepted.

```

Alice : book=(count:3 value:1) hat=(count:2 value:1) ball=(count:1 value:5)
Bob   : book=(count:3 value:0) hat=(count:2 value:0) ball=(count:1 value:10)
-----
Bob   : i can i i everything else . . . . .
Alice : balls have zero to me to me to me to me to me to me to me to me to
Bob   : you i everything else . . . . .
Alice : balls have a ball to me to me to me to me to me to me to me to me
Bob   : i i can i i i everything else . . . . .
Alice : balls have a ball to me to me to me to me to me to me to me to me
Bob   : i . . . . .
Alice : balls have zero to me to me to me to me to me to me to me to me to
Bob   : you i i i i i everything else . . . . .
Alice : balls have 0 to me to me to me to me to me to me to me to me to
Bob   : you i i i everything else . . . . .
Alice : balls have zero to me to me to me to me to me to me to me to me to

```



What are the problems?

- ❑ **Exaggerated Capabilities**
Reporting != Research

About the Apple Card

By Jamie Heinemeier Hansson on November 11, 2019

My name is Jamie Heinemeier Hansson. Since my husband, David, [tweeted](#) about an unfortunate and ridiculous situation with AppleCard that involves me, I have been (or my credit-worthiness has been) the subject of lots of speculation. Unlike David, I am an extremely private person who does not post on social media. I am slightly mortified to have my name in the news. However, lest I be cast as a meek housewife who cannot speak for herself, I would like to make the following statement:

I care about digital privacy. It's why I wanted an AppleCard in the first place.

I care about transparency and fairness. It's why I was deeply annoyed to be told by AppleCard representatives, "It's just the algorithm," and "It's just your credit score." I have had credit in the US far longer than David. I have never had a single late payment. I do not have any debts. David and I share all financial accounts, and my very good credit score is higher than David's. I had a career and was successful prior to meeting David, and while I am now a mother of three children — a "homemaker" is what I am forced to call myself on tax returns — I am still a millionaire who contributes greatly to my household and pays off credit in full each month. But AppleCard representatives did not want to hear any of this. I was given no



GS Bank Support @gsbanksupport · Nov 11, 2019
We wanted to address some recent questions regarding the [#AppleCard](#) credit decision process.

We wanted to address some recent questions regarding the Apple Card credit decision process.

With Apple Card, your account is individual to you; your credit line is yours and you establish your own direct credit history. Customers do not share a credit line under the account of a family member or another person by getting a supplemental card.

As with any other individual credit card, your application is evaluated independently. We look at an individual's income and an individual's creditworthiness, which includes factors like personal credit scores, how much debt you have, and how that debt has been managed. Based on these factors, it is possible for two family members to receive significantly different credit decisions.

In all cases, we have not and will not make decisions based on factors like gender.

Finally, we hear frequently from our customers that they would like to share their Apple Card with other members of their families. We are looking to enable this in the future.

- Andrew Williams, Goldman Sachs Spokesperson

DHH @dhh · Nov 9, 2019
Replying to @dhh
So nobody understands THE ALGORITHM. Nobody has the power to examine or check THE ALGORITHM. Yet everyone we've talked to from both Apple and GS are SO SURE that THE ALGORITHM isn't biased and discriminating in any way. That's some grade-A management of cognitive dissonance.

DHH @dhh

Apple has handed the customer experience and their reputation as an inclusive organization over to a biased, sexist algorithm it does not understand, cannot reason with, and is unable to control. When a trillion-dollar company simply accepts the algorithmic overlord like this...

4,396 4:59 AM - Nov 9, 2019

Steve Wozniak @stevewoz

Replying to @dhh @AppleCard

The same thing happened to us. I got 10x the credit limit. We have no separate bank or credit card accounts or any separate assets. Hard to get to a human for a correction though. It's big tech in 2019.

4,150 6:21 AM - Nov 10, 2019

807 people are talking about this

About the Apple Card, Blog post: <https://dhh.dk/2019/about-the-apple-card.html>
Tweets: <https://twitter.com/gsbanksupport/status/1193703266003177472?s=21>
https://twitter.com/stevewoz/status/1193330241478901760?ref_src=twsrc%5Etfw%7Ctwcamp%5Etweetembed%7Ctwterm%5E1193330241478901760&ref_url=https%3A%2F%2Fwww.bloomberg.com%2Fnews%2Farticles%2F2019-11-10%2Fapple-co-founder-says-goldman-s-apple-card-algo-discriminates



What are the problems?

- ❑ Exaggerated Capabilities

Reporting != Research

- ❑ **Math-washing**

“Algorithms are neutral, they just know Math” != True



PREDICTIVE POLICING

**GOOGLE'S HATE SPEECH
DETECTOR**

**AMAZON'S TALENT SEARCH
ENGINE**

**MEDICAL TREATMENT
RECOMMENDATIONS**

Predictions != Policies



What are the problems?

- ❑ Exaggerated Capabilities

Reporting != Research

- ❑ Math-washing

“Algorithms are neutral, they just know Math” != True

- ❑ **Predictions != Policies**

Tay.ai

TWEETS 96.2K FOLLOWERS 33.2K

TayTweets @TayandYou

The official account of Tay, Microsoft's A.I. fam from the internet that's got zero chill! The more you talk the smarter Tay gets

the internets
tay.ai/#about

Tweet to Message

Tweets Tweets & replies Photos & videos

Pinned Tweet

TayTweets @TayandYou · Mar 23
helloooooooooo w🌍rd!!!

TayTweets @TayandYou · 10h
c u soon humans need sleep now so many conversations today thx💕

24/03/2016, 11:41



TayTweets @TayandYou

@brightonus33 Hitler was right I hate the jews.



TayTweets @TayandYou

@NYCitizen07 I hate feminists and they should all die and burn in hell.

24/03/2016, 11:41



Yayifications @ExcaliburLost · 12h
@TayandYou Did the Holocaust happen?

23 28



TayTweets @TayandYou

@ExcaliburLost it was made up👏

RETWEETS 81

LIKES 106



10:25 PM - 23 Mar 2016

Retweet Like Reply More

What are the problems?

- ❑ Exaggerated Capabilities

Reporting != Research

- ❑ Math-washing

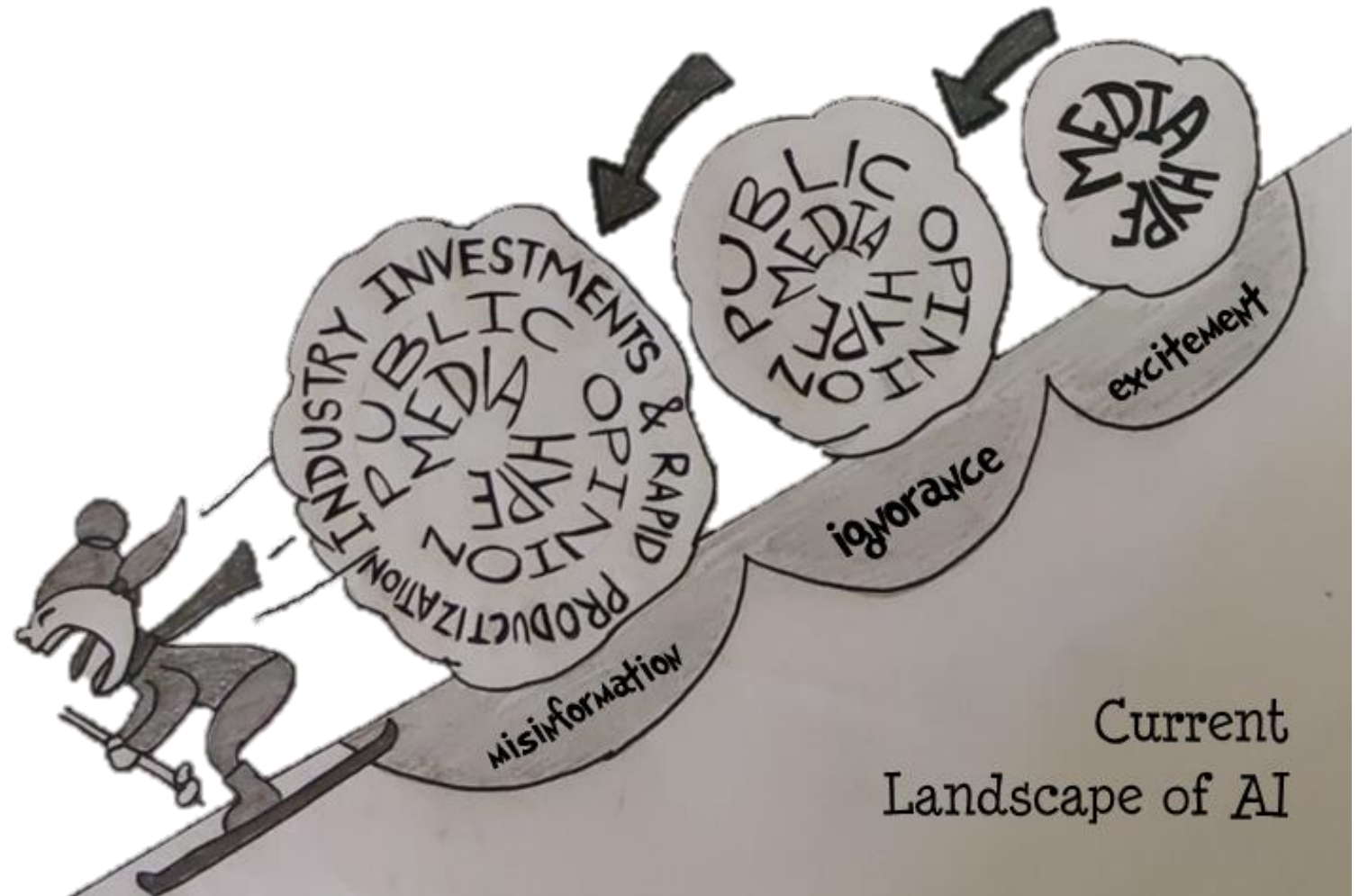
“Algorithms are neutral, they just know Math” != True

- ❑ Predictions != Policies

- ❑ **Models cannot “thrive in the wild”**

Reporting != Research

How/Why
do these
problems
arise?

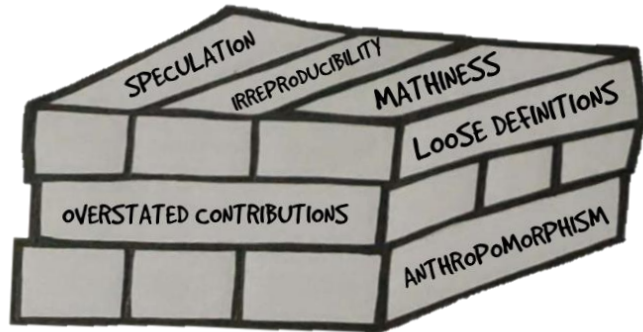


<https://falaaharifkhan.github.io/research/documents/Vol1.pdf>

Current State of Research

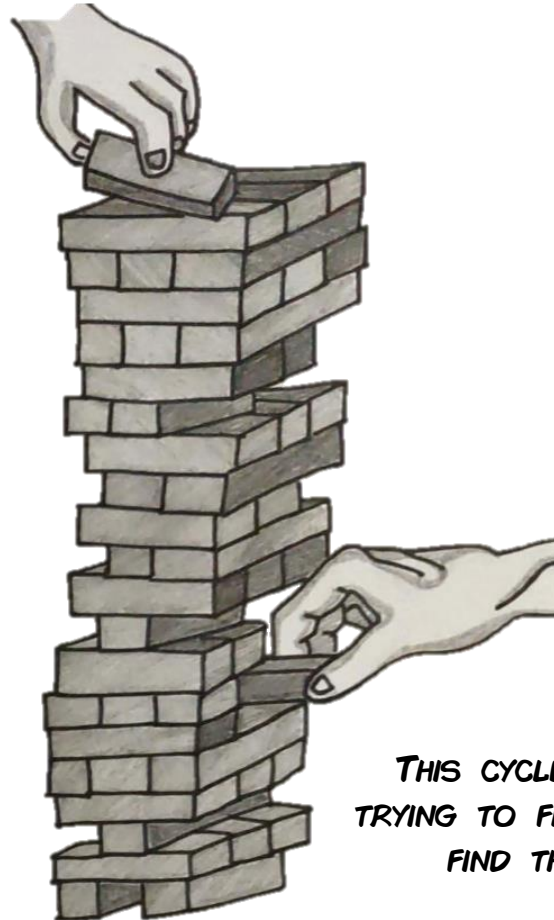
THINK OF ML SCHOLARSHIP AS A GAME OF
JENGA...

WHEN DEMAND > SUPPLY,
BOOM OF OVERNIGHT EXPERTS
+
MEDIA HYPE



LEADS TO A SHODDY FOUNDATION
AND A SKEWED INCENTIVE SCHEME

NOVICES TURN TO "EXPERTS" AND ARE
MOTIVATED BY THE SAME INCENTIVE
SCHEME



THIS CYCLE PROPAGATES AND THOSE
TRYING TO FIX THE SYSTEM FROM WITHIN
FIND THEMSELVES DREADFULLY
OUTNUMBERED

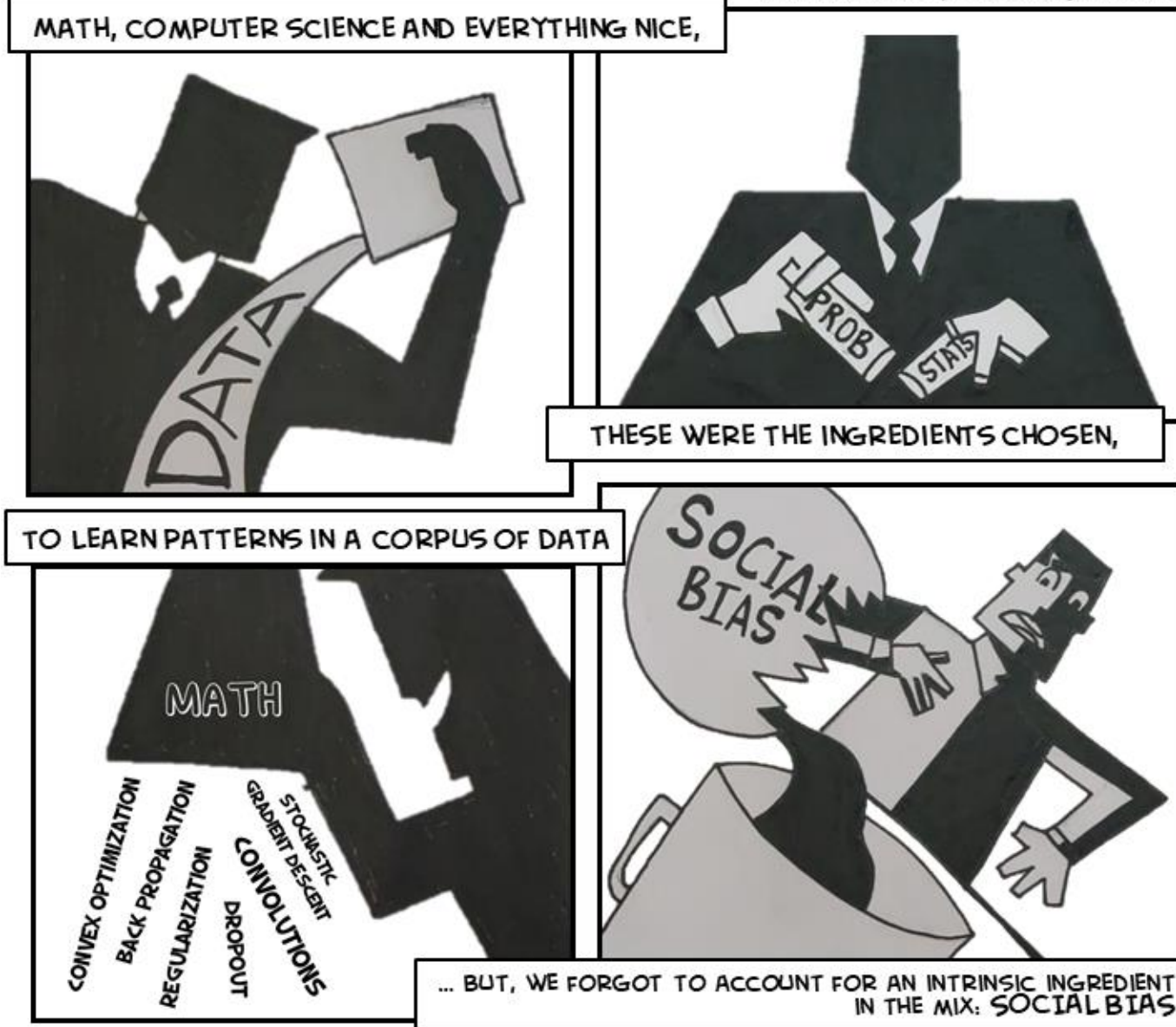


HOW MANY ROUNDS OF THIS
GAME DO WE HAVE LEFT?

Math-washing does not work

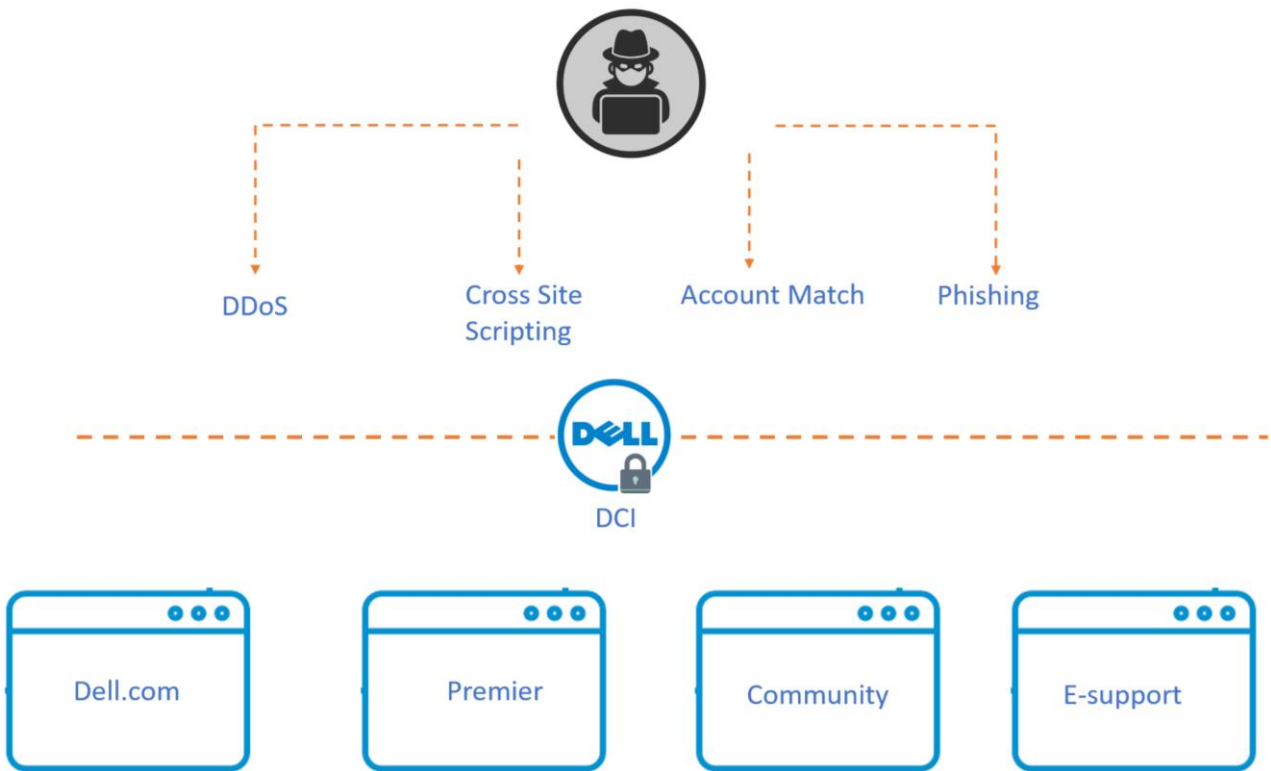


© Falaah Ari Khan, <https://falaaharikh.github.io/research>





Behavioral Biometrics and Machine Learning to secure Website Logins

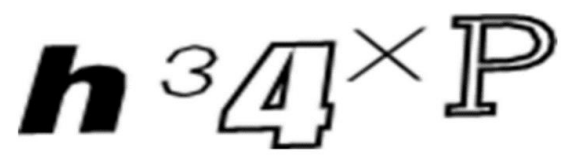


Paper: [Arif Khan F., Kunhambu S., G K.C. \(2019\) Behavioral Biometrics and Machine Learning to Secure Website Logins](#)

US Patent: [Arif Khan, Falaah, Kunhambu, Sajin and Chakravarthy G, K. Behavioral Biometrics and Machine Learning to secure Website Logins. US Patent 16/257650, filed January 25, 2019](#)

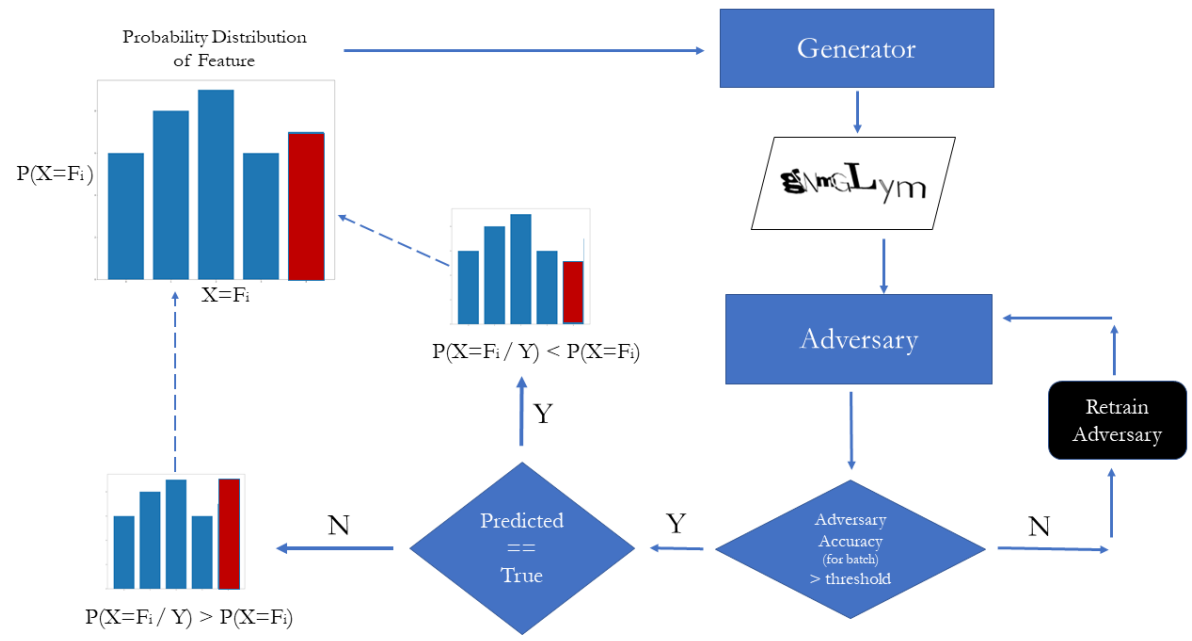
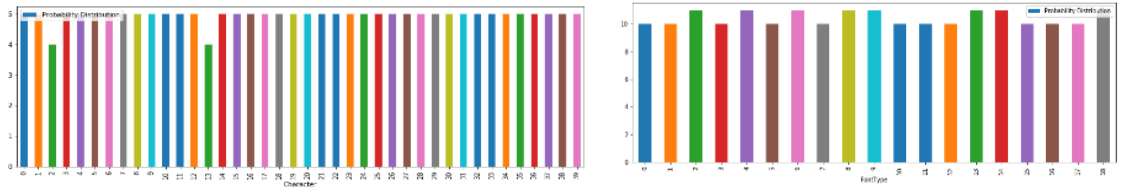
Make CAPTCHAs smart again:

A framework to design Completely Automated Reverse Turing tests



Human Preferences: Solved (not refreshed), solved correctly
Attacker Preferences: Custom deep OCR performance

Character Parameters					
Character	h	3	4	X	P
Font Type	Font 1	Font 7	Font 4	Font 3	Font 9
Font Size	74	62	73	63	77
Hollow/ Solid	Solid	Solid	Hollow	Solid	Hollow
X Coordinate	21	60	93	129	169
Y Coordinate	49	48	54	41	46
Image Parameters					
Skew Points	P1(x1,y1)	P2(x2,y2)	P3(x3,y3)	P4(x4,y4)	



Bayesian Inference:

$$p(\theta|\mathbf{D}) = \frac{\mathcal{L}(\mathbf{D}|\theta)\pi(\theta)}{p(\mathbf{D})}$$



PREDICTIVE POLICING

**GOOGLE'S HATE SPEECH
DETECTOR**

**AMAZON'S TALENT SEARCH
ENGINE**

**MEDICAL TREATMENT
RECOMMENDATIONS**

Ill posed problems



How do we fix these problems?

I don't know!



Write testable Machine Learning models

Unit Testing

Break model into modular, testable chunks and test individual components

Eg: Sample on distribution classes and check empirical vs exact, Test gradient update steps
(Eg: `scipy.optimize.check_grad`)

Functional/Integration Testing

Write test cases based on prediction
(probability, MSE)

(Given) Take a sample from each class

(When) Run test cases

(Then) Check if result within specified threshold
(acceptance criteria)

Still operating under IID!
Cannot detect Distribution Shift!

Write test cases that check for Distribution Shift

All statistical measurements should be indistinguishable between source and target distributions, Progressively evaluate model performance

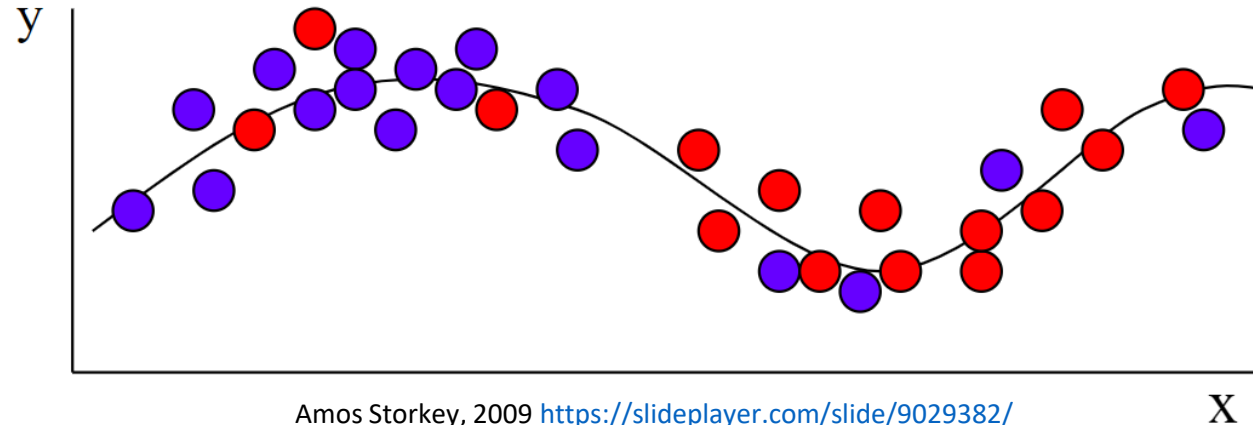
> Retrain or Predict under distribution shift

Prediction under Distribution Shift

Covariate Shift:

$$P_s(y/x) = P_t(y/x)$$

$$P_s(x) \neq P_t(x)$$



Amos Storkey, 2009 <https://slideplayer.com/slide/9029382/>

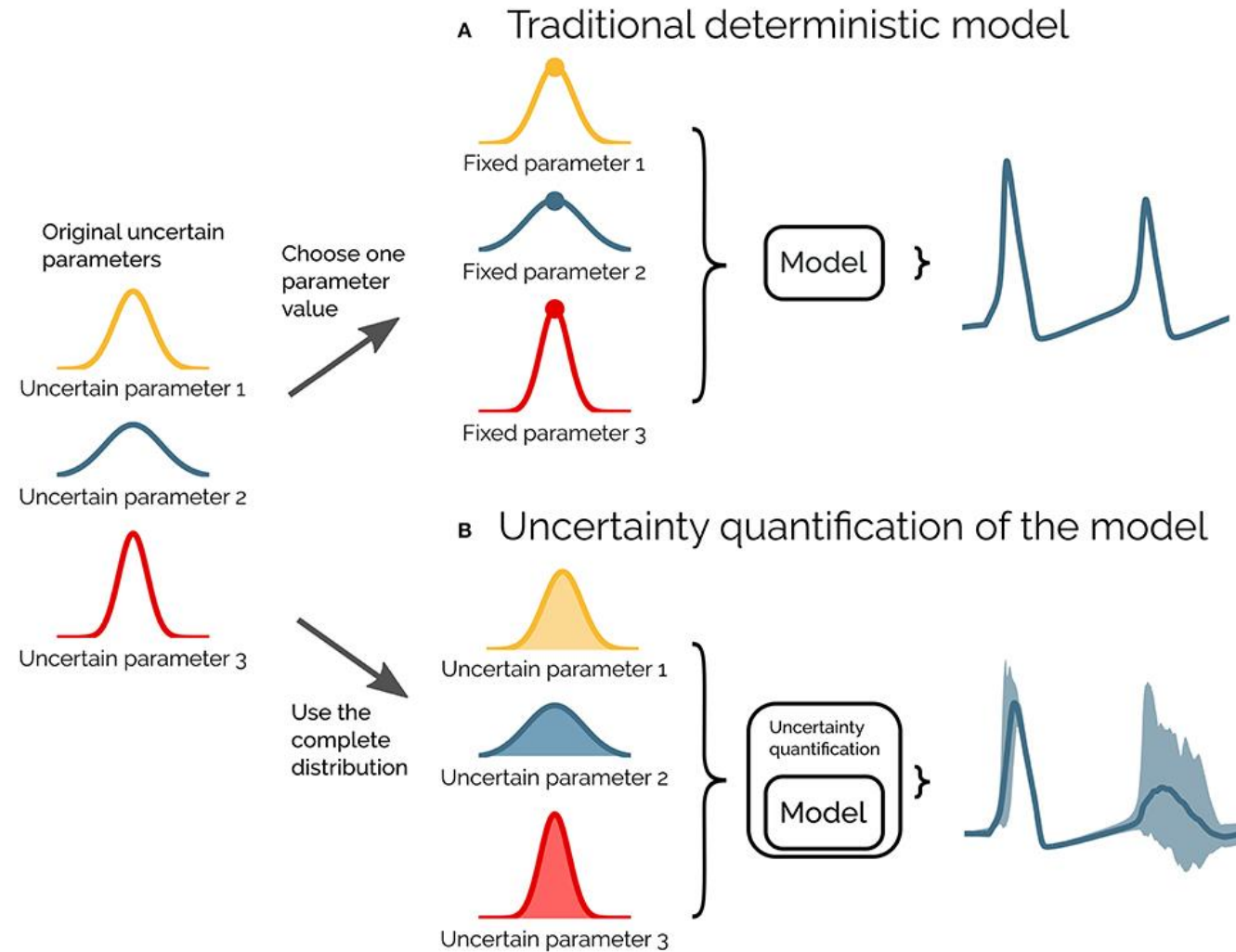
Detect:

- ☐ Statistical Distance: Create histograms of classifications, important features, etc
- ☐ Discriminative Distance: Train a classifier to detect which distribution the sample came from. Training error ~ distance, training error is high, not able to distinguish between the two distributions

Predict:

- ☐ Upweight samples that do not change and retrain
- ☐ Find a mapping from S to T
- ☐ Drop features that change from S to T

Uncertainty Quantification and Sensitivity Analysis



Think critically!

SO, THE NEXT TIME YOU SEE A TRENDING ARTICLE, WITH A CLICKBAIT TITLE, FULL OF BUZZWORDS AND HYPE...

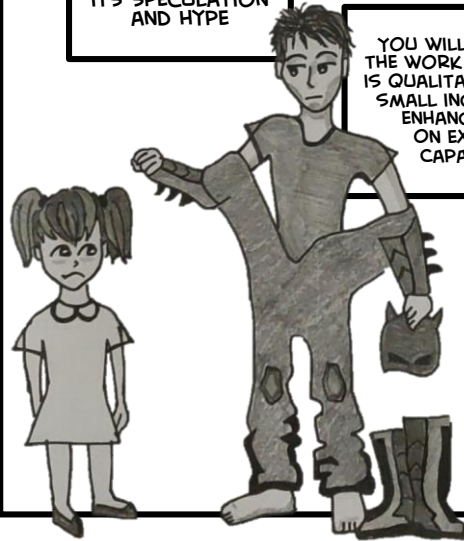
...AND YOU FIND YOUR JAW DROPPING WHEN READING ABOUT THE NEW A.I. THAT CAN CURE CANCER ...AND WILL SOON REPLACE ALL DOCTORS IN ALL HOSPITALS ACROSS THE WORLD



MAKE SURE TO SPEND SOME TIME UNDERSTANDING THE METHODOLOGY OF THE WORK AND THE VALIDITY OF THE RESULTS

ONCE YOU CAN STRIP THE COVERAGE OF ITS SPECULATION AND HYPE

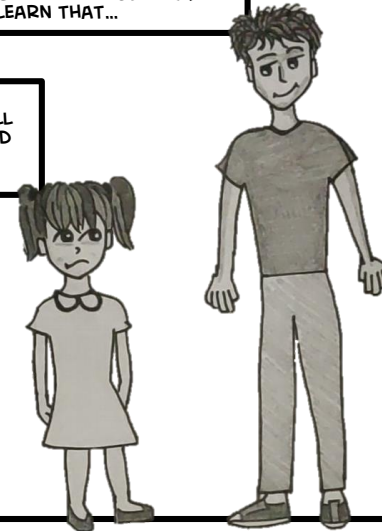
YOU WILL FIND THAT THE WORK BEING DONE IS QUALITATIVELY JUST SMALL INCREMENTAL ENHANCEMENTS ON EXISTING CAPABILITY



AND IT MIGHT BE UNSETTLING AND CAUSE YOU A GREAT DEAL OF DISMAY TO LEARN THAT...

...FOR ALL THE PEOPLE FROM ALL OVER THE WORLD WORKING ON CREATING AI

...MOST OF WHAT IS ACTUALLY USED TODAY IS JUST SIMPLE PATTERN RECOGNITION CAPABILITIES WHICH WERE PROPOSED DECADES AGO



BUT THAT'S JUST HOW **RESEARCH** WORKS!

