

### Data Science – how it should look

Huge amounts of novel social data

Variety, velocity, volume!

ML tools – automate analysis

Death of theory



### But – in practice...

Usually not VVV.

These methods really work best for operational development and process modelling

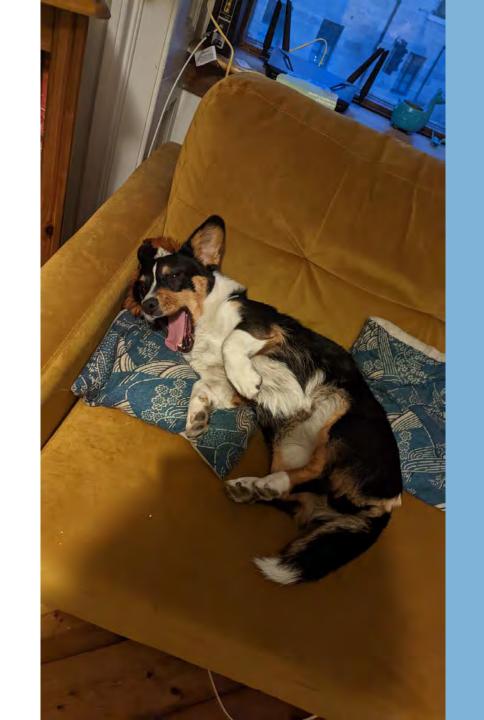
Insights – less so

A lot of the point of data science is feature reduction

- the errors cancel out

Theory's back, and it's good again

Much more rubbish, interesting, and complicated than we thought



### Cambridge Cybercrime Centre

Set up in Cambridge University Computer
Lab

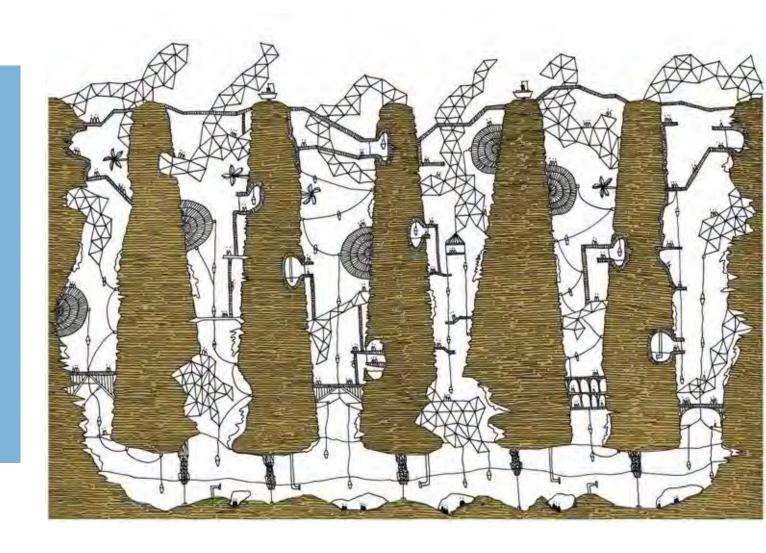
Problem – lots of scrappy data sources

Collect huge amounts of 'authoritative' data and license to researchers

Do our own research

Try to implement the Cathy O'Neill approach (old fashioned now!)

My attempts to do bad data science wel



What are we going to do with all this these our data?

Look at the cybercrime-as-a-service ecosystem

Huge transformation about 15 years ago

Cybercrime now about small businesses running services

What effects do law enforcement interventions have and

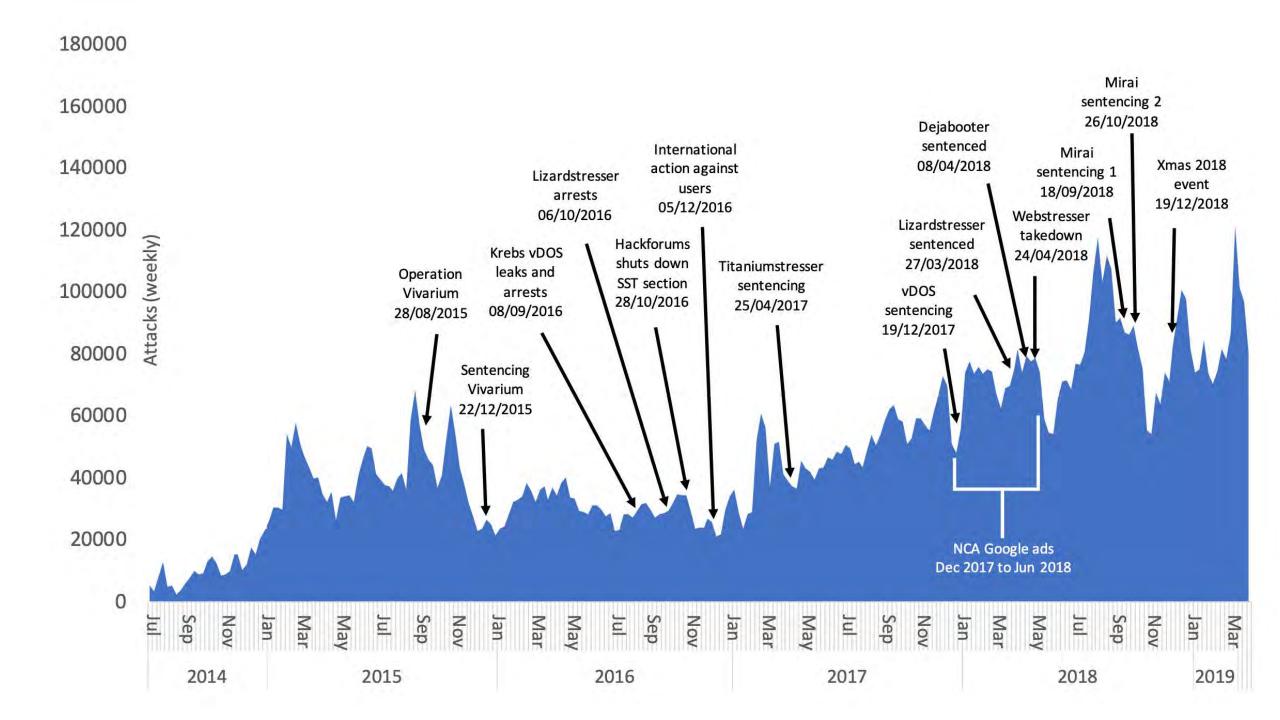
why?

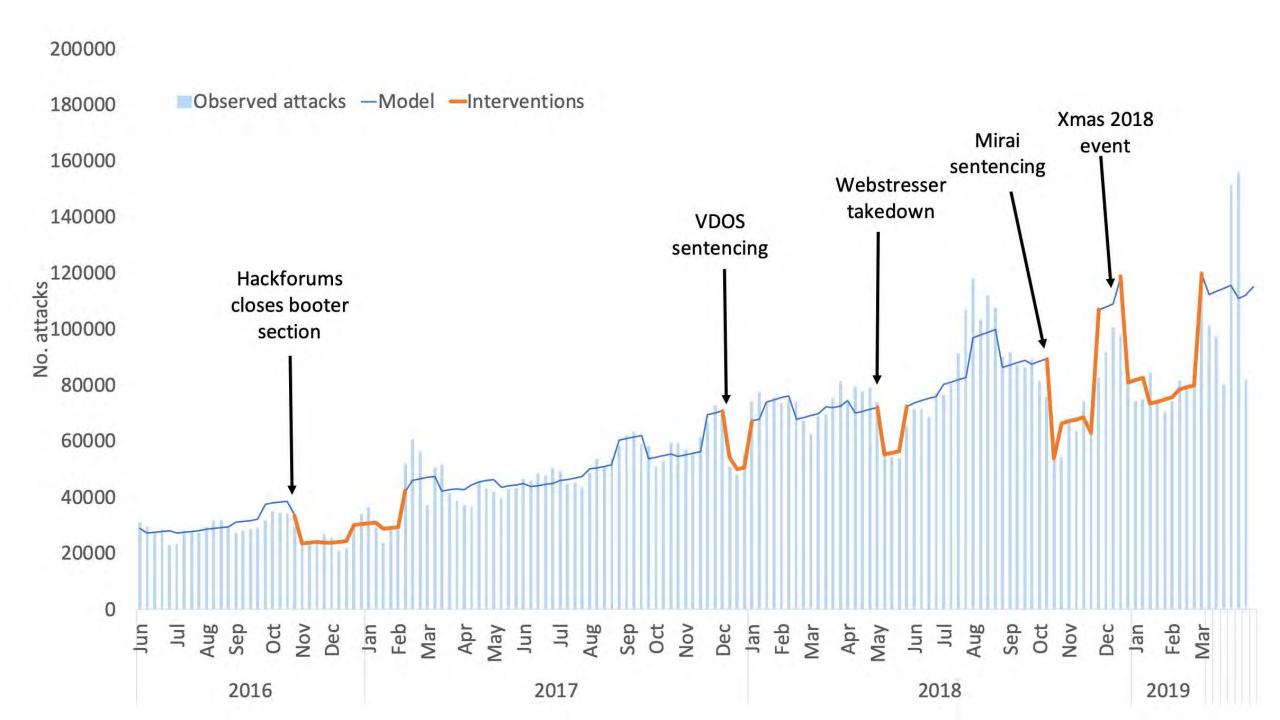


## Honeypots – measuring crime directly

Set up servers which pretend to be crime infrastructure

Linear time series – count of attacks per country





Intervention		UK	US	RU	FR	DE	PL	NL	Overall
$\begin{array}{c} Xmas 2018 \\ Intervention \\ 19/12/2018 \end{array}$	Mean L95/U95 Duration Signif.	-27% -43/-28% 9 weeks 0.000**	-49% -55/-42% 9 weeks 0.000**	-33% -43/-22% 9 weeks 0.000**	-1% -13/11% N/A 0.828	-28% -36/-20% 8 weeks 0.000**	-23% -37/-5% 3 weeks 0.014*	-16% -27/-3% 8 weeks 0.018*	-32% -37/-27% 10 weeks 0.000**
Mirai sentencing and other actions 24/10/2018	Mean L95/U95 Duration Signif.	-27% -42/-9% 2 weeks 0.006**	-31% -41-20% 7 weeks 0.000**	-5% -16/7% 2 weeks 0.41	-9% -31/21% N/A 0.533	-32% -40/-23% 6 weeks 0.000**	-47% -56/-36% 2 weeks 0.000**	-19% -35/0% 6 weeks 0.053	-40% -46/-34% 8 weeks 0.000**
Webstresser takedown 24/04/2018	Mean L95/U95 Duration Signif.	-10% -21%/3% N/A 0.120	-24% -40/-4% 4 weeks 0.022*	-16% -33/6% 2 weeks 0.14	-22% -35/-7% 4 weeks 0.006*	-29% -36/-22% 9 weeks 0.000**	-29% -42/-14% 6 weeks 0.001**	146% 94/211% 4 weeks 0.000**	-21% -30/-12% 3 weeks 0.000**
vDOS sentencing 16/12/2017	Mean L95/U95 Duration Signif.	-20% -33/-5% 3 weeks 0.011*	-4% -18/12% 3 weeks 0.563	-37% -47/-24% 2 weeks 0.000**	-30% -37/-23% 2 weeks 0.000**	-4% -17/10% N/A 0.532	16% -17/62% N/A 0.373	-24% -33/-13% 3 weeks 0.000*	-24% -32/-25% 3 weeks 0.000**
HackForums 28/10/2016	Mean L95/U95 Duration Signif.	-48% -53/-42% 15 weeks 0.000**	-30% -37/-21% 7 weeks 0.000**	-13% -23/-3% 14 weeks 0.02*	-52% -59/-43% 15 weeks 0.000**	-32% -41/-23% 7 weeks 0.000*	2% -19/28% N/A 0.86	-35% -42/-27% 15 weeks 0.000*	-30% -33/-25% 13 weeks 0.000**

and the second control of the second control

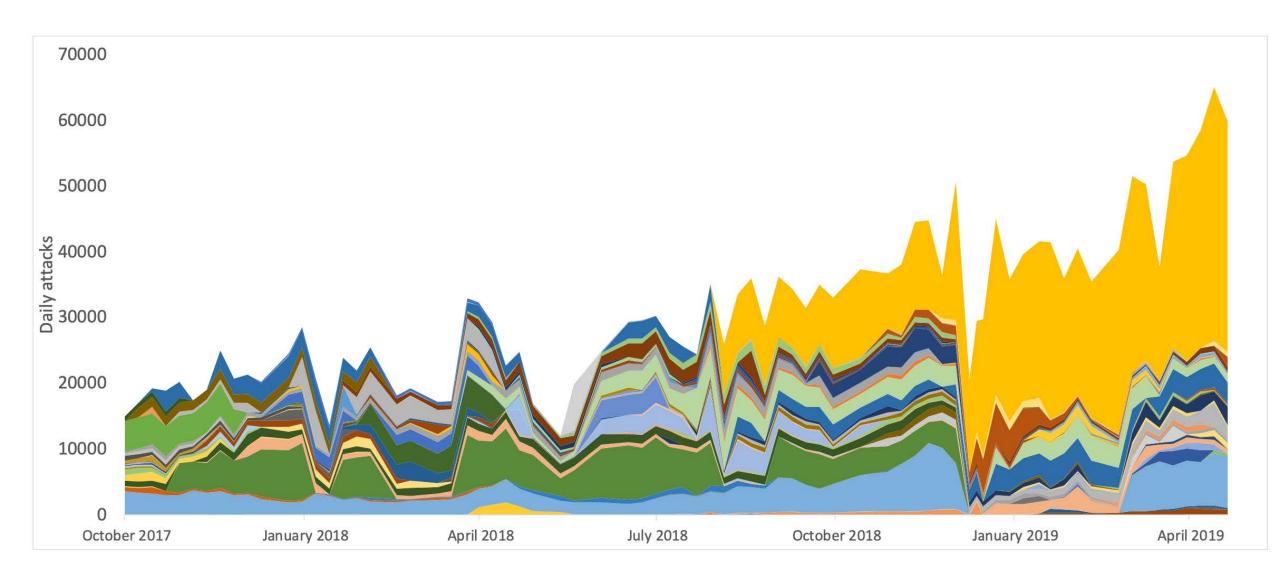
## Secret bonus analysis

We had also been collecting what the illegal services *said* they were doing – they advertise a weekly number of attacks

Taking notes on a criminal conspiracy...

Journal hated it, but we used some fun stats magic to 'prove' that it was 'reliable' and they relented

Normally distributed and heteroskedastic



Begin initial analysis – problem generation

Started by trying to expain our quants findings

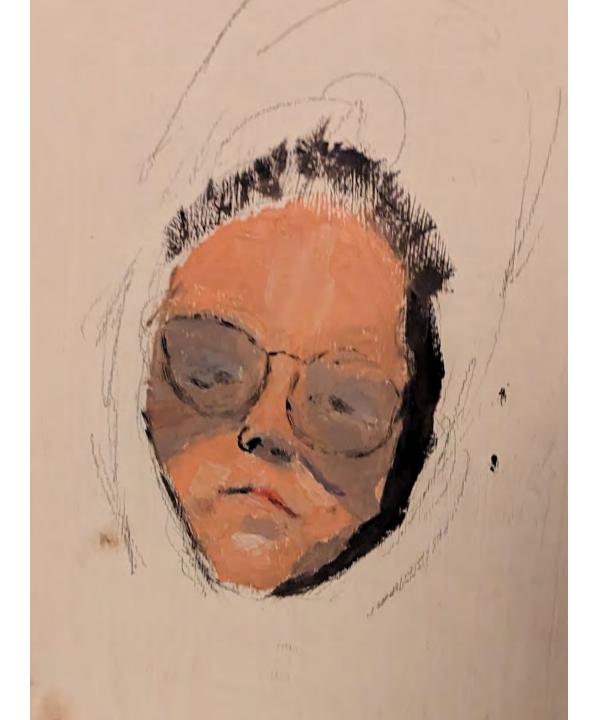
Why do some of these strategies work and some don't?

Hugely overdetermined – we can see what works (ish) but no idea

why

Could write a rubbish quants paper at this point but decide not to

Carry out more research!



### Interviews

Interview people in the communities running these attack platforms and users Easier than you think – but you get lied to a lot



## CrimeBB – data science

Flagship collection of cybercrime forums
Home-made scrapers

Millions of posts across 10 years and 30 forums



## CrimeBB – enrichment

Sentiment
Post-type
Crime type
Topics
Group based trajectory modelling
Data tools

### In practice...

Research often just doesn't work

#### **QUANTS BAD**

Our quants collections were pretty good, but incredibly overdetermined.

#### **QUAL BAD**

Our qual collection was based on quite a small number of interviews...

### DATA SCIENCE WORSE

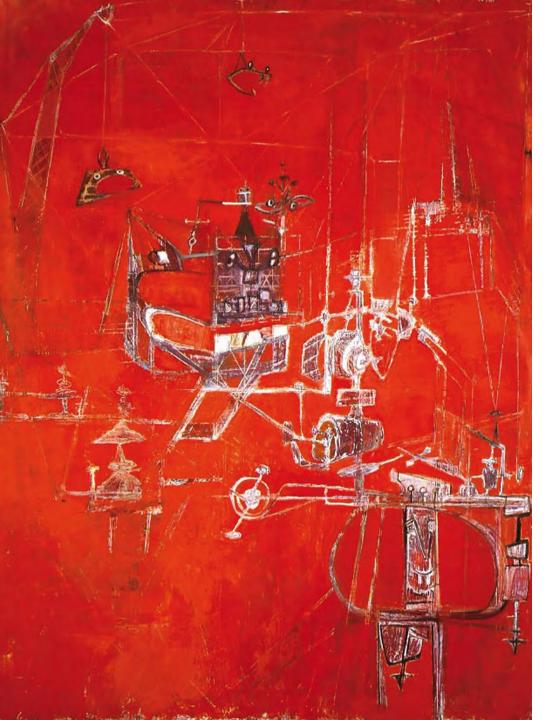
And our data science modelling was a nightmare. NLP fell over just looking at the complexity and size of our data, and our GBTM models first refused to converge, then worked but wrecked my hard drive, and finally gave results that look like every other GBTM in the world...

# Social Data Science to the rescue!

Take what was good about our collections – the range of different types of knowledge

Iterate back and forward,



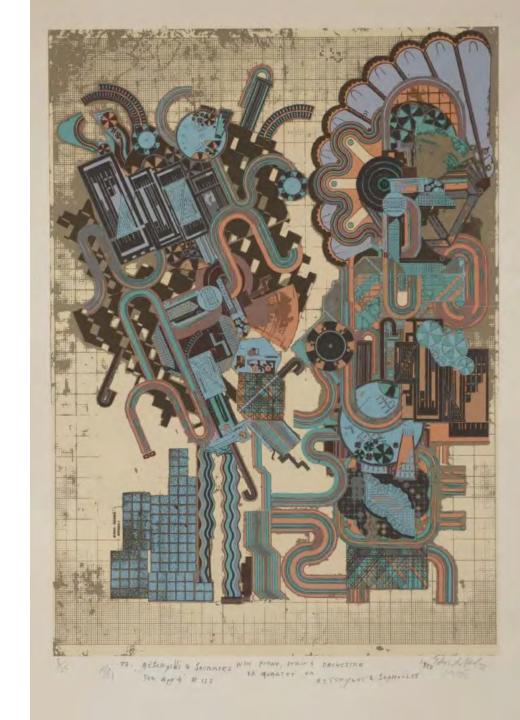


### First – we immerse

- Read the forums manually every day for 20 minutes and a good scan of the new chats
- Try to mimic a lurker
- Tool design round 1 has been done
- But need to build tools to probe the dataset
- Think what do I need? What can't I get just from looking?

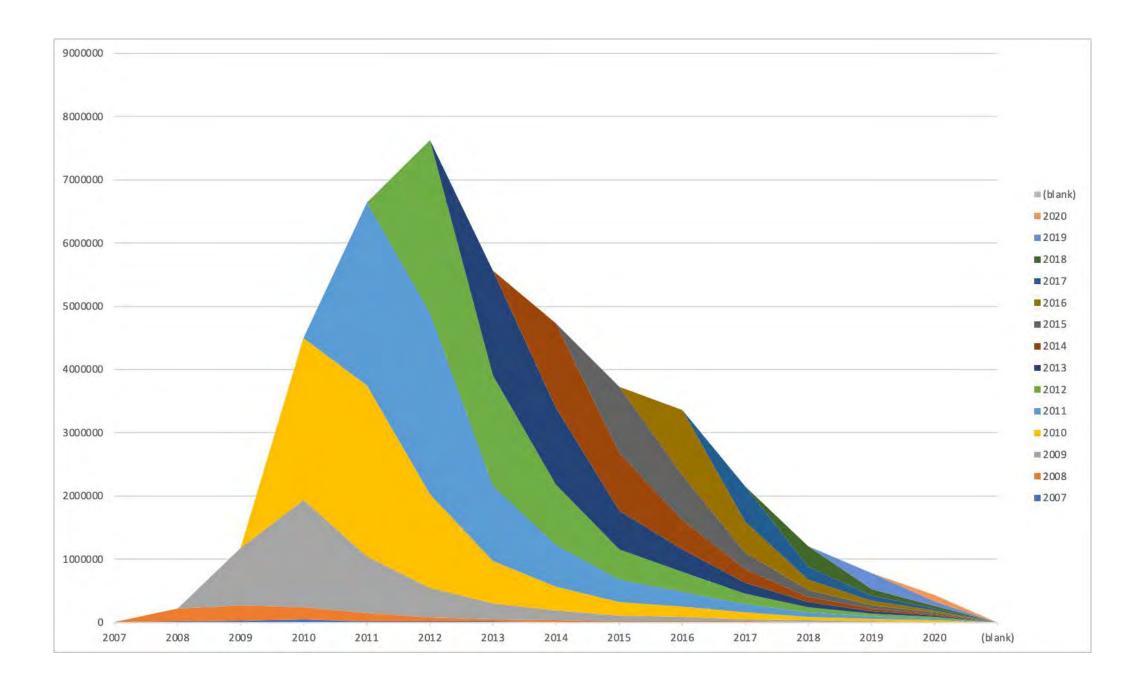
### Second – we iterate

- Working with both kinds of data as a team
- Weekly meetings ('stand ups') and over-the-shoulder chats
- Draw out the evolving 'models' of the codes, concepts, and situations on shared whiteboards
- Bring interim findings and codes back and forward to one another
- Eventually these stabilize each model should have a list of assertions/findings and the linked evidence
- Plan two papers (one technical and one sociological) from the start - sieve findings into buckets then group buckets within papers



### Second – we iterate

- Even active people don't stay in the community for more than a year except in exceptional cases – and permanence getting lower as time goes on
- Huge range of topics hacker culture, lessons, marketplace etc. Lots of cybercrime services
- Drill down into a particular subforum, then subset again for keywords like 'work' (or which reference particular kinds of work)
- Scan through and sample posts, find language used, then run queries and see if it looks like what you thought it would
- Get a subset of 'meaty' (relevant) discussions
- Apply to wider dataset manually (or topic model within your subset, then apply the topic classifier to the broader dataset)
- Group based trajectory modelling (more or less a time-series LCA)





# Third: initial interviews

Map practices, people, actors, technologies, ideas

Idea of boredom comes up a lot

Hacker ideas still really important – skill etc.

But actual work is very simple



### Stage 4: deep-dive

- Expand the map
- Take the concepts from the interviews and search them
- Theory-driven (Star, labour alienation)
- How do key words change over time?
- How is involvement changing over time?
- Can we search for posts using the language we're seeing in the interviews?
- Sort topic models then sample within topics, then model topics within topics etc.
- Can we model posts like this versus other sorts of posts?
- Establishing our rough key themes
  - Very few posts talk about genuinely technical stuff mostly about low-level, boring, or facile work
  - Lots of evidence of boredom
  - Work is mostly maintenance and administration, or marketing and customer service
  - Technical talk getting less over time
  - Evidence of deskilling and gradual breakdown

### Go back to interviews

Stage 5: exploration

Pull findings from the data – established some key themes now

Explore qualitative depth, new findings from data

## Findings

People are bored – directly	contradicts	the core	literature,	which	views	hacking
as exciting						

Mostly administrative work

No respect

Exit through burnout (you can even look at people's last posts)

None of our sources was authoritative on its own – but we could evidence the wider applicability of our interview findings using the quantitative and data science data

Why do some police actions work and not others? Making the work even more boring! Alienation of labour.

Why does cybercrime increase during the pandemic, but not really change its nature? Well-developed service market

# Immersive data ethnography

AKA 'social data science"

Focus on multiple types of data

Immersion – how was data produced

Patch up the rubbish by attacking your data from all directions

READ THE MICRODATA

A lot of this is about how you write it up

## THANK YOU

@johnnyhistone