A DATA SCIENCE ANABASIS



THIS IS A STORY ABOUT A REAL LIFE DATA SCIENCE **PROJECT.**

- DATA SCIENCE PROJECTS ARE WEIRD.
- THEY DON'T FIT WITH HOW COMPANIES WORK.
- BUT OCCASIONALLY SOME REAL GEMS FALL OUT **UNEXPECTEDLY**.
 - THIS IS A STORY ABOUT ONE OF THOSE GEMS.
 - AND THE JOURNEY THAT TOOK US THERE.

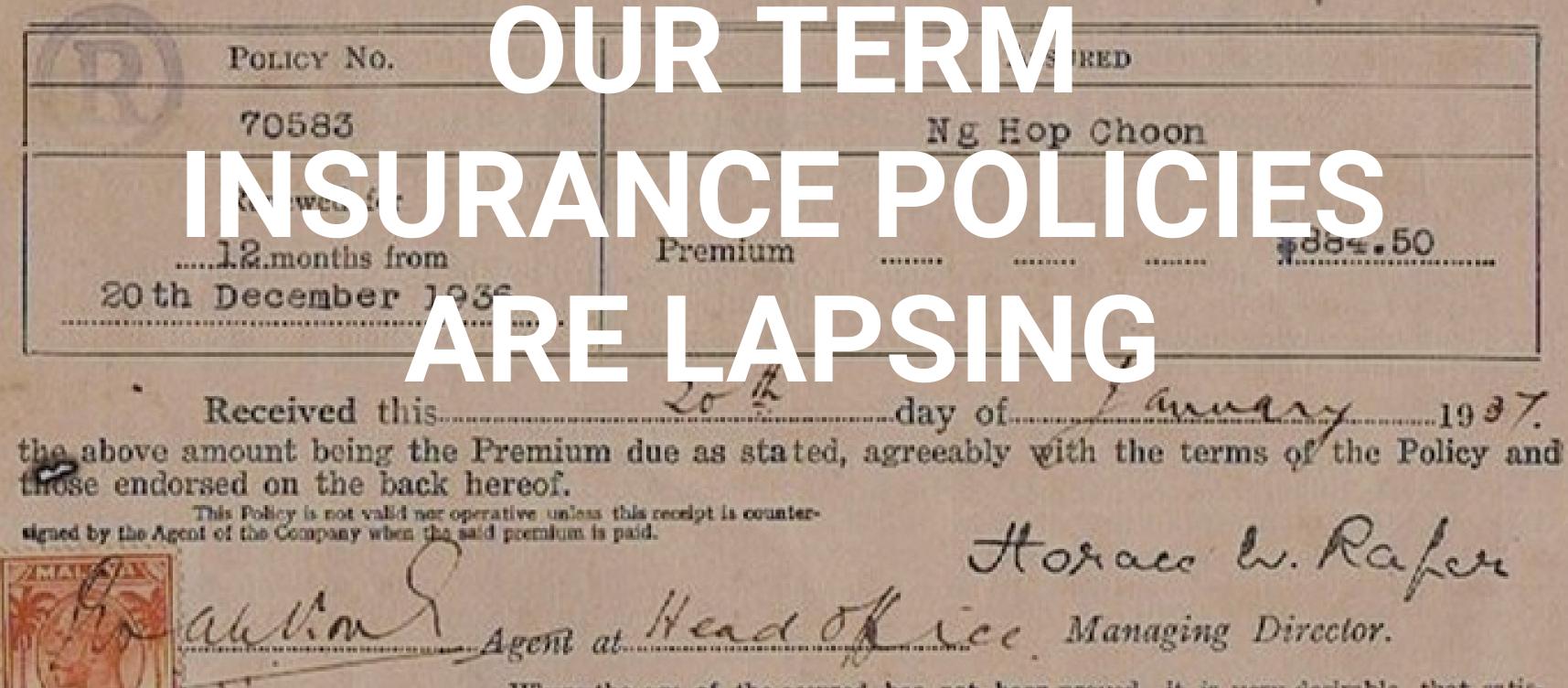
THE PROBLEM





The Great Eastern Like Assurance Co., Ltd.

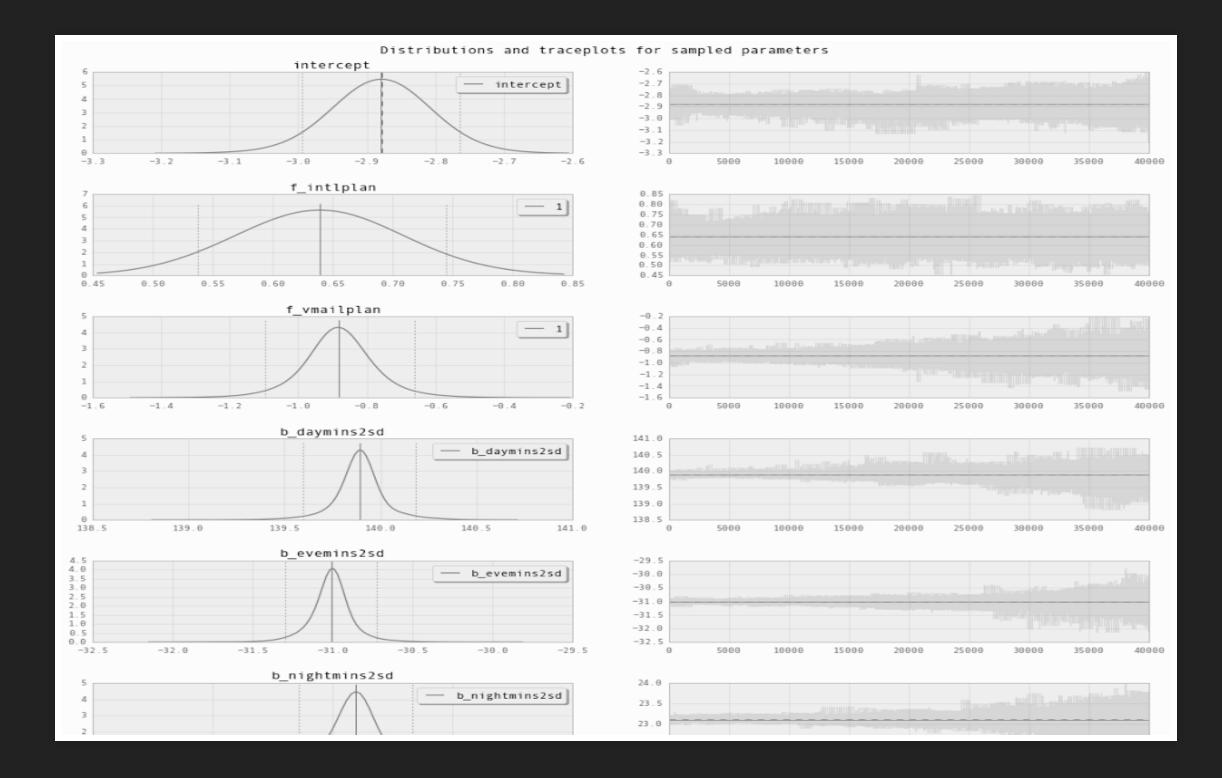
(Incorporated in the Straits Settlements.) HEAD OFFICE: - - SINGAPORE.

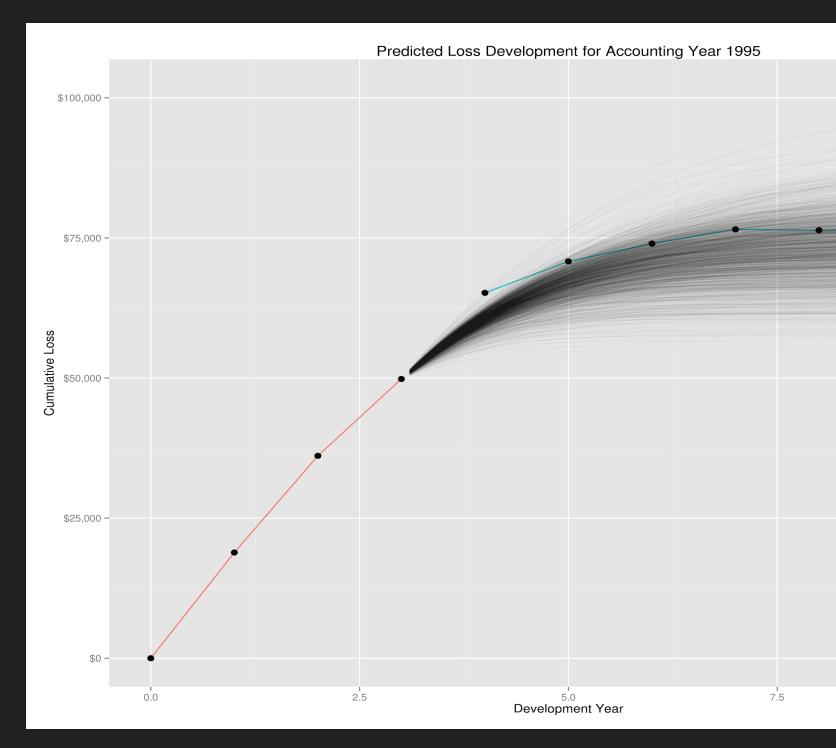


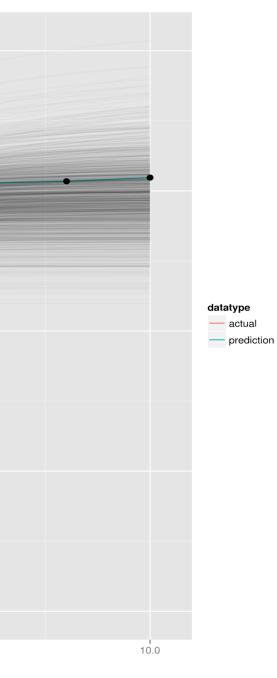
Where the age of the assured has not been proved, it is very desirable that satisfactory evidence of age be produced to the Company, so that the age may be admitted during

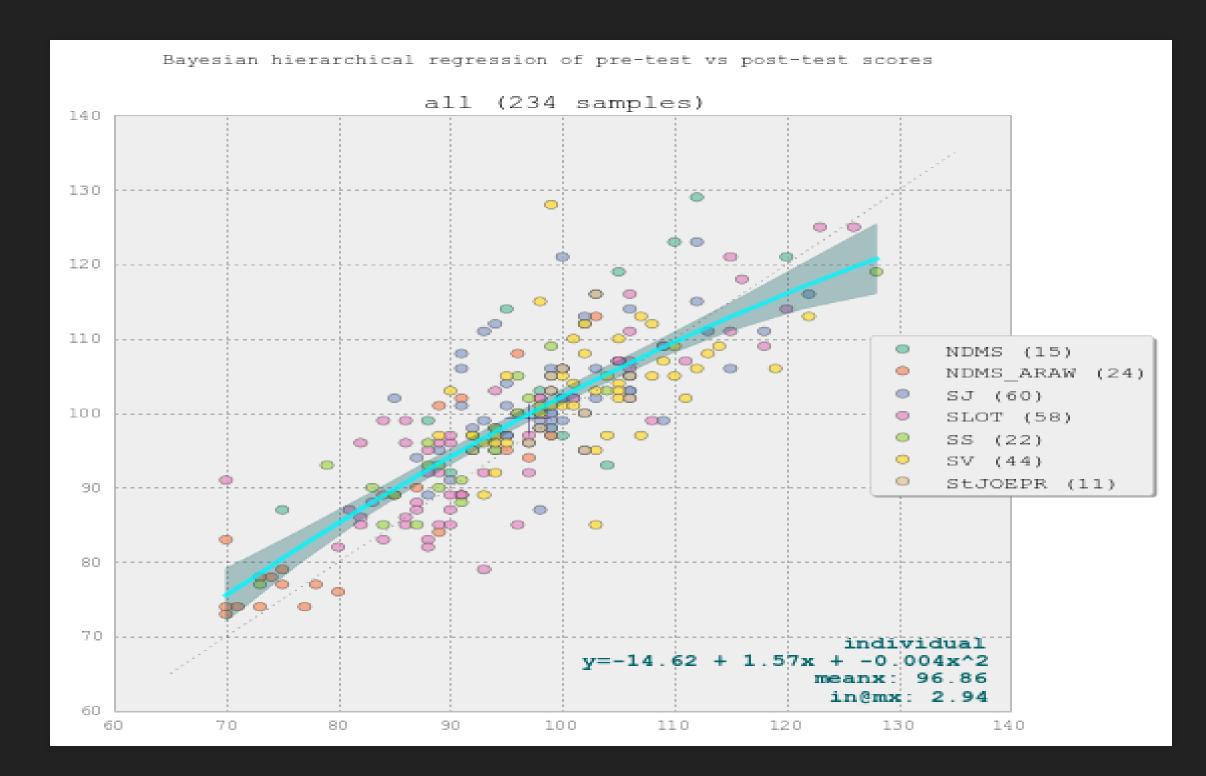
THE BRIEF

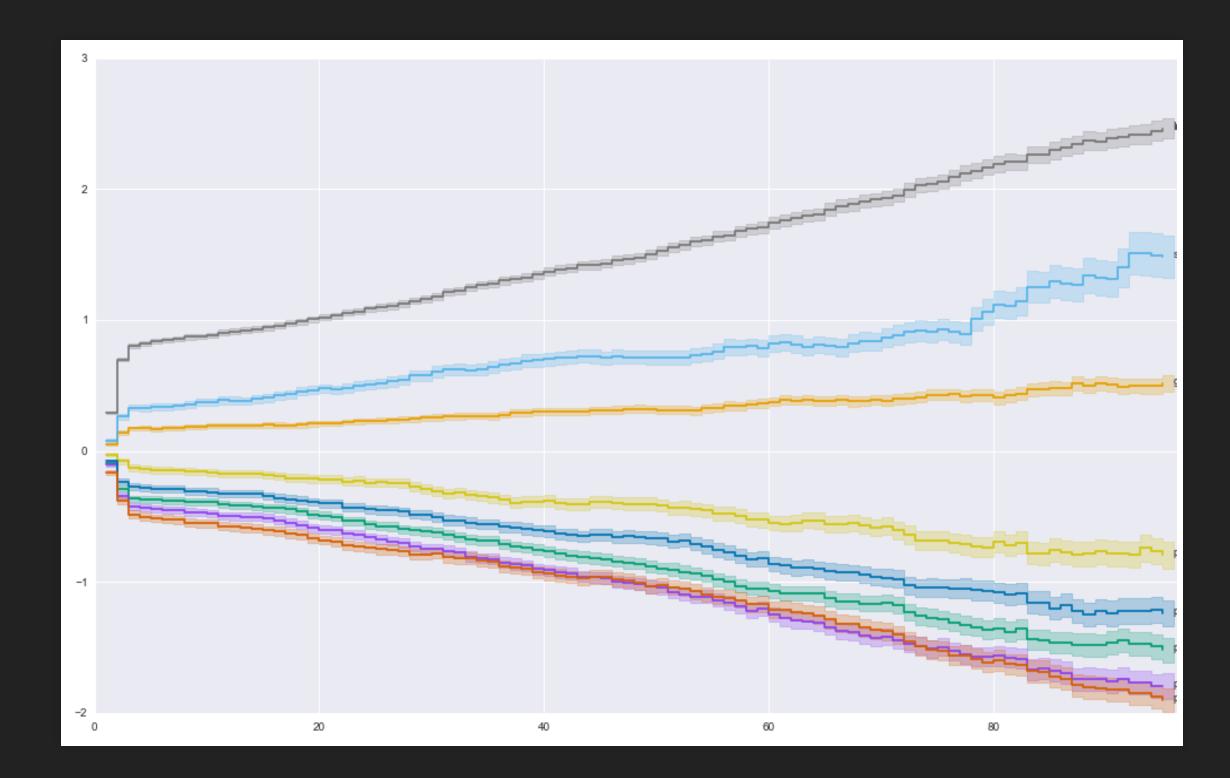
- Client: We want you to investigate the causes of lapse for a few hundred thousand policies.
- Us: Great sounds like a cool project!
- Client: We want to use a Bayesian approach cause we think its cool!...
- ...and we met this cool guy at a conference he invented MCMC...
- ...and he wants a holiday in Ireland and you get to work with him - cool!.
- Us: Cool what could possibly go wrong!











elcocustomerchurn ×											
3888/notebooks	/demo_telcod	ustom	erchurn.ipynb	#							
IP[y]:	Note	bod	ok de	emo_telc	ocustomerc	hurn					
File Edit View Insert Cell Kernel Help											
Code + Cell Toolbar: None +											
Out[113]:		state	accountdur	areacode	phonenumber	intlplan	vmailplan	vmailmsg	daym		
	customerid										
	408-327- 4579	GA	111	408	327-4579	1	0	0	159.4		
	408-327- 4579	GA	111	408	327-4579	1	0	0	159.4		
	408-327- 4579	GA	111	408	327-4579	1	0	0	159.4		
	408-327- 4579	GA	111	408	327-4579	1	0	0	159.4		
	408-327- 4579	GA	111	408	327-4579	1	0	0	159.4		

5 rows × 36 columns

mins	daycalls	daycharge	 evecharge2
.4	47	27.1	 0.674905
.4	47	27.1	 0.674905
.4	47	27.1	 0.674905
.4	47	27.1	 0.674905
.4	47	27.1	 0.674905

A COUPLE OF MONTHS LATER



AWKWARD QUESTIONS



WHAT HAPPENED?

- The Bayesian approach just didn't work.
- We were overawed by academic fire power...
- ...and spent too long flogging a dead horse.

Lesson 1: When you're in a hole stop digging...

...trust your judgement and go back to the...



ON A POSITIVE NOTE

- We learned a lot about...
 - The problem
 - The data
 - Bayesian statistics always useful
 - The inherent computational complexities
 Ninja level Python

Lesson 2: Try and avoid doing computer science in the office

useful complexities

WHAT WE DID NEXT



WHILE WAS THERE

I looked up a really smart lady I'd met at R in

Insurance.. ...who runs the Risk Center at University of Barcelona. She told me that our approach would not work and... .gave me a paper on Time Varying Survival Analysis I thought what's the worst that can happen?

OPSUCCESS

- The



THIS ALREADY SUFFICIENT TO ANSWER

WHAT IS THE EXPECTED RANGE OF SURVIVAL OVER TIME?

CAN I PREDICT THE SURVIVAL RATE? WHAT ARE THE DRIVERS OF THE SURVIVAL RATE

RESULTS OF SURVIVAL ANALYSIS

- We now had a good grasp on what makes policies lapse
- But the company had limited information on their clients
- Things like:
 - Age, Sex, Smoker
 - Term & Premium
 - Policy Type & Broker
 - Premium payment history

THE COMPANY HAD A CREDIT **RISK PROBLEM**

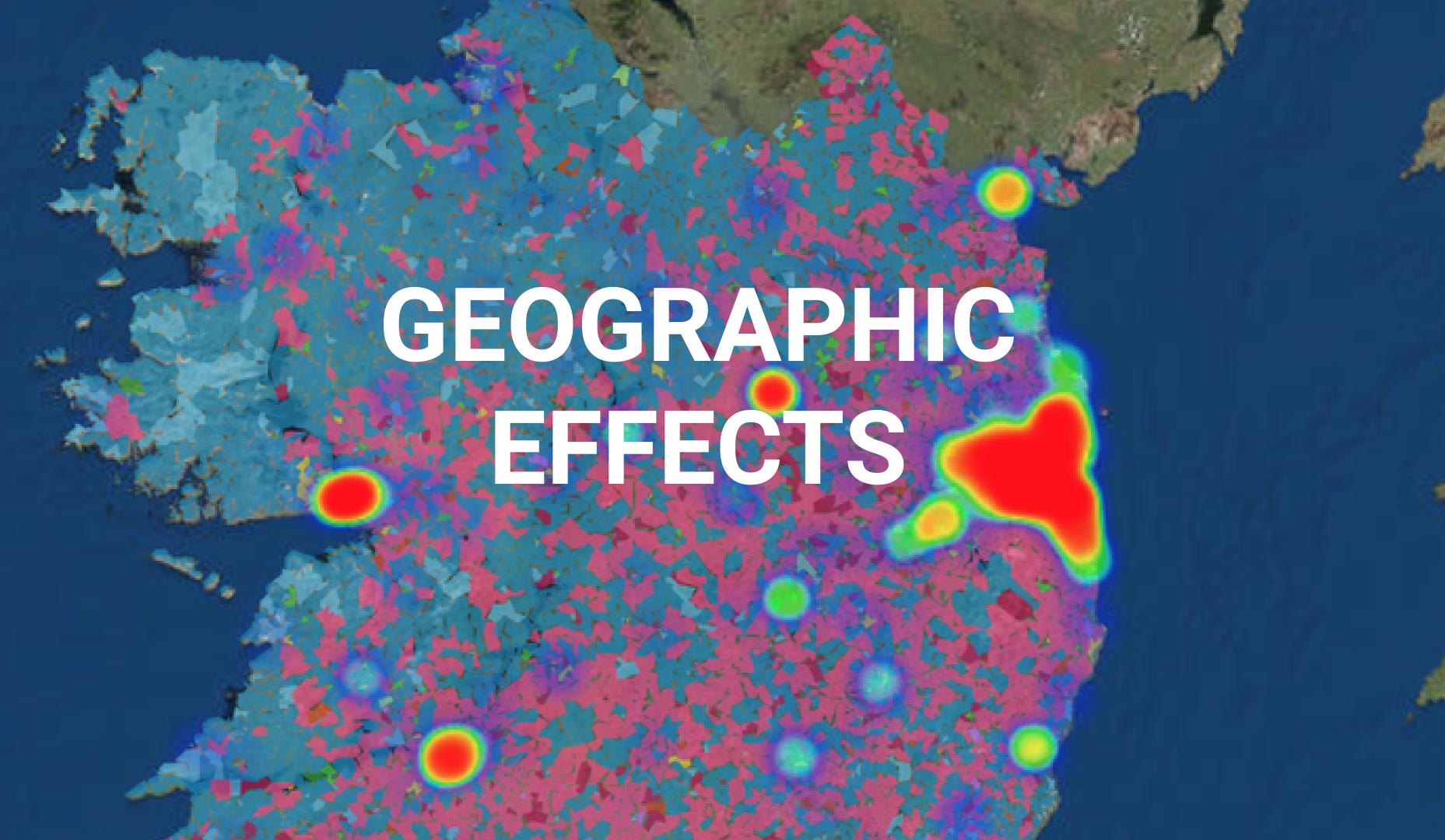
WHAT DID WE LEARN?

- Talk to everybody and filter later
- I can't emphasise this next point enough...
- ...go to...
 - Meetups
 - Technical User Groups
 - Conferences
- You'll learn a lot from other people...
- ...and every now and then a random conversation will save your life

ARE WE DONE YET? ALONG THE WAY WE HAD NOTICED A FEW THINGS







SOCIOEGONOMIC EFFECTS



HOW CAN WE USE THEM? • TO ENCOURAGE CUSTOMERS TO STAY • TO HELP PRICE RISK • IDENTIFY NEW MARKETS SO WE CALLED EXPERIAN

IT WAS A SHORT CONVERSATION

HI, I'D LIKE SOCIOECONOMIC INFORMATION FOR 250K **ADDRESSES**

HOW MUCH!

OKAY SO THAT'S A NON-RUNNER



GEOCODING

• FIRST LET'S GEO-CODE OUR ADDRESSES

• WE HAD TWO CHOICES:

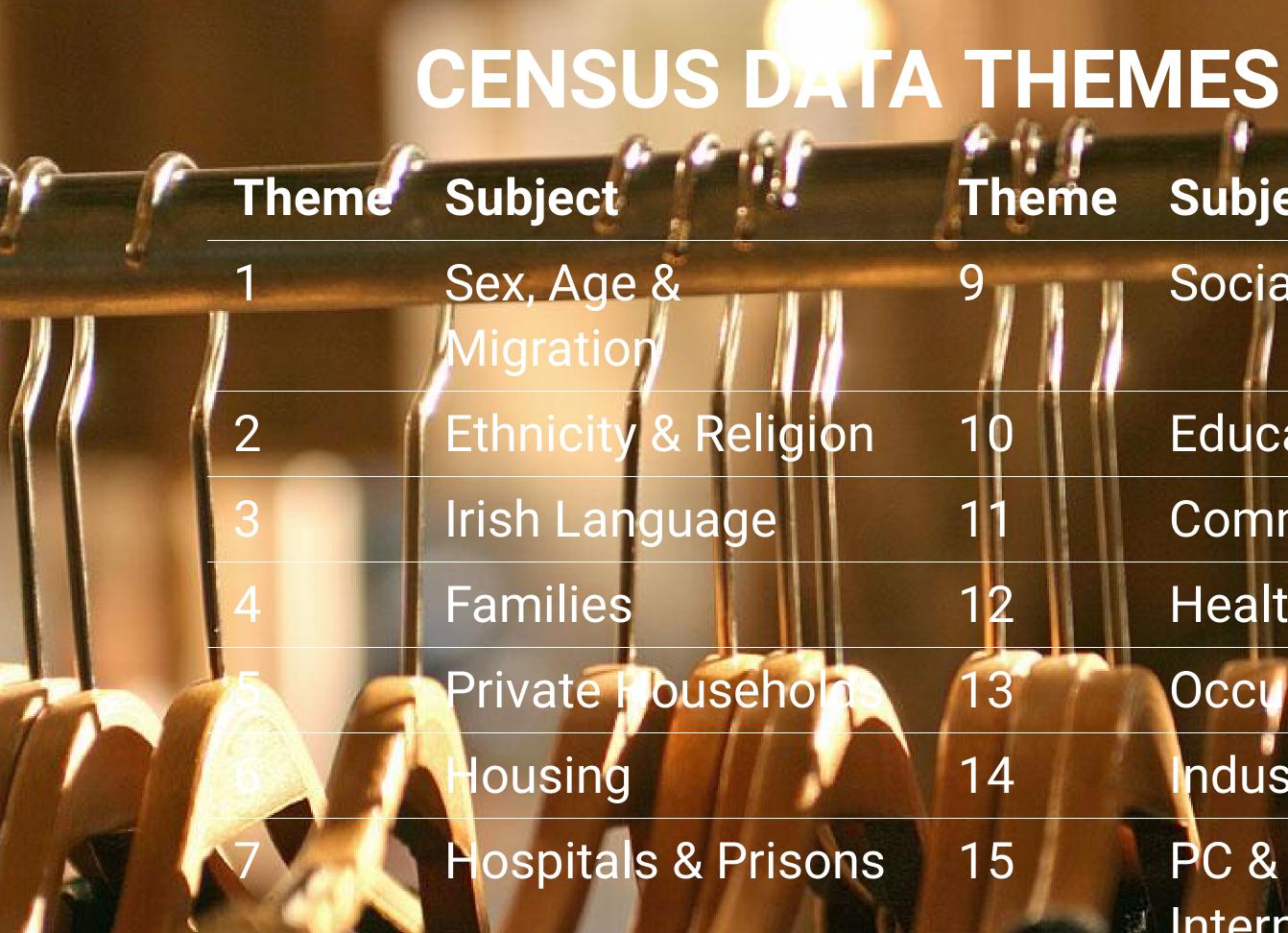
Use Google - which you pay for Use Nominatim - FOSS / roll your own

WE TRIED BOTH FOR IRELAND, GOOGLE IS BETTER MAINLY BECAUSE

IRISH ADDRESSES ARE PATHOLOGICAL! BUT WE NOW HAVE A LAT / LONG FOR EACH CLIENT

SHOPPING FOR DATA THERE'S A LOT OF IT AROUND WE LOOKED AT THE CENSUS DATA IN IRELAND A CENSUS IS DONE EVERY FIVE YEARS THE AMOUNT OF INFORMATION IN IT IS ASTOUNDING **APPROXIMATELY 700 FEATURES COVERING 15 THEMES**

WHAT DATA IS AVAILABLE



Subject Social class Education Commuting Health Occupation ndustries PC & Internet



SMALL AREA MAPS

THIS IS TERRIEYING

THE SMALLEST OUTPUT AREA FOR CENSUS DATA • ~20,000 SMALL AREAS COVERING IRELAND • EACH COVERS APPROX. 200 PEOPLE • EACH CENSUS FEATURE AVAILABLE AT THIS LEVEL

WE REALSED THAT

WE COULD DO WHAT EXPERIAN DOES WE WOULD HAVE THE CODE • WE COULD INTEGRATE IT WITH ANY DATA SCIENCE PROJECT

• WE COULD TUNE IT TO FIT OUR PARTICULAR NEEDS

LET'S GO



NOTI A START OF A CONTRACTOR A C

IT'S HARD TO MAKE SENSE OF THIS MUCH DATA:

THERE ARE 18,488 SMALL AREA MAPS

EACH SMALL AREA MAP IS REPRESENTED BY A ROW

EACH ROW HAS 767 ENTRIES ONE FOR EACH FEATURE

b_1 a_{1n} a_{11} a_{12} EAS, WE EAVE b_2 a_2 BG WITH 18,488 ROWS & 767 COLUMNS

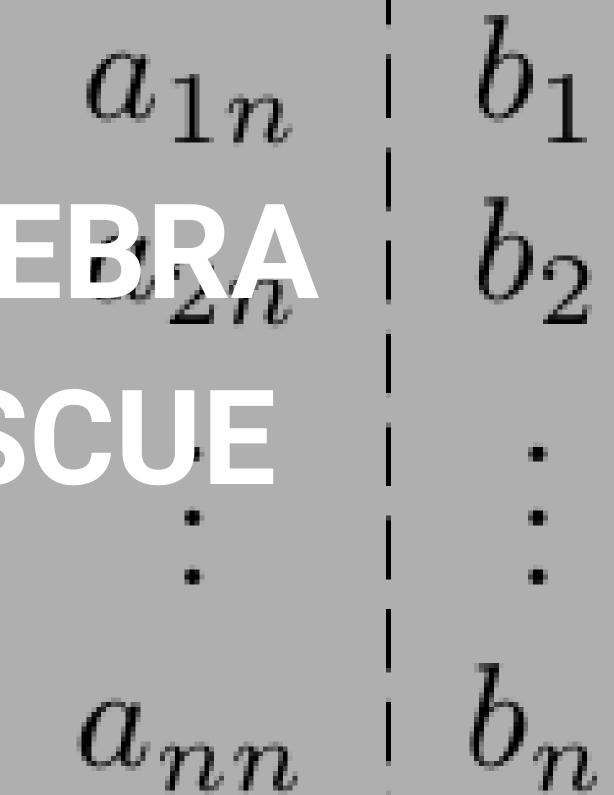
 a_{n1} a_{n2}

a_{nn}

b_n

a_{11} a_{12} LABAR-ALGEBRA a_{21} TOTHE RESCUE

 a_{n1} a_{n2}

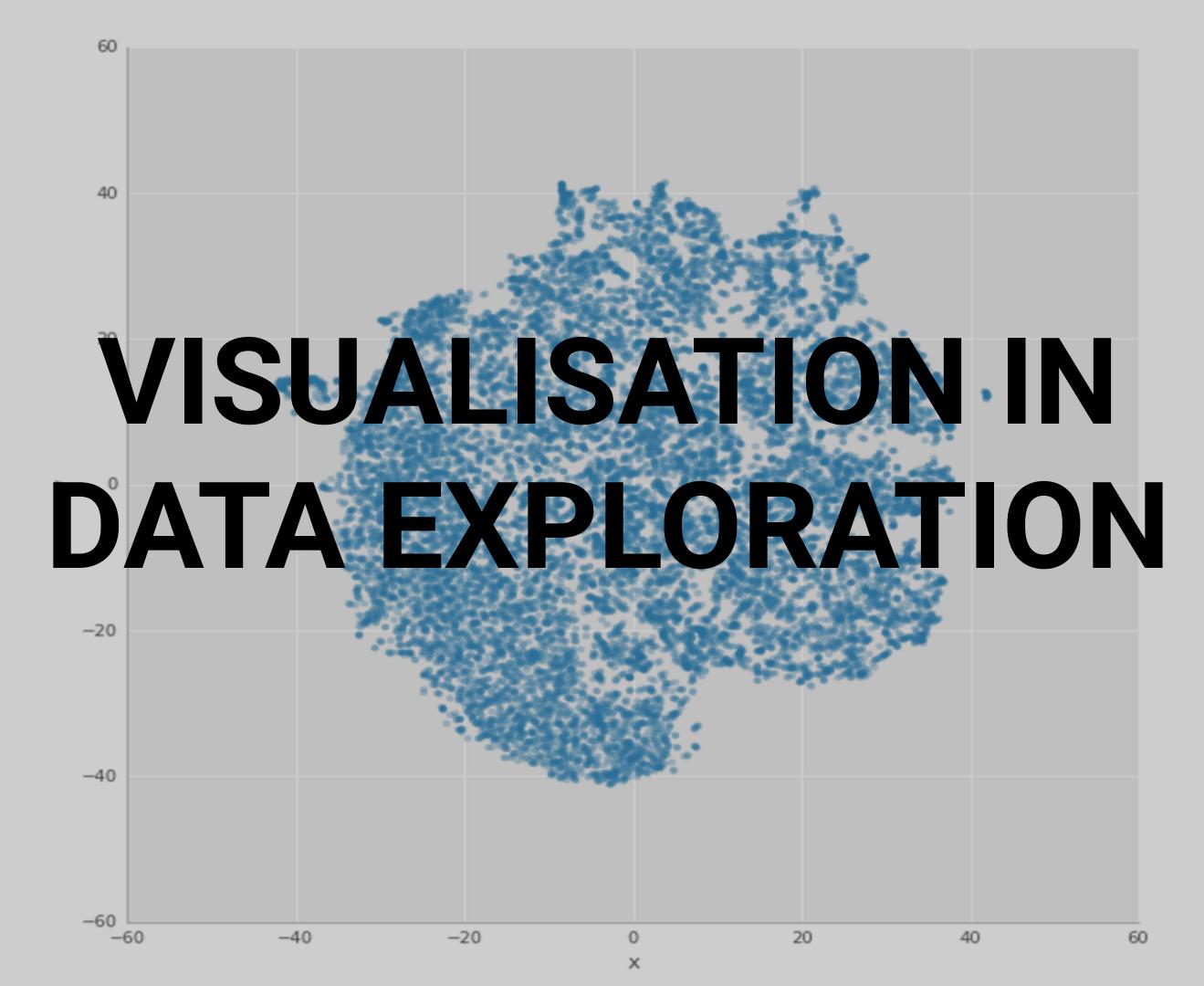


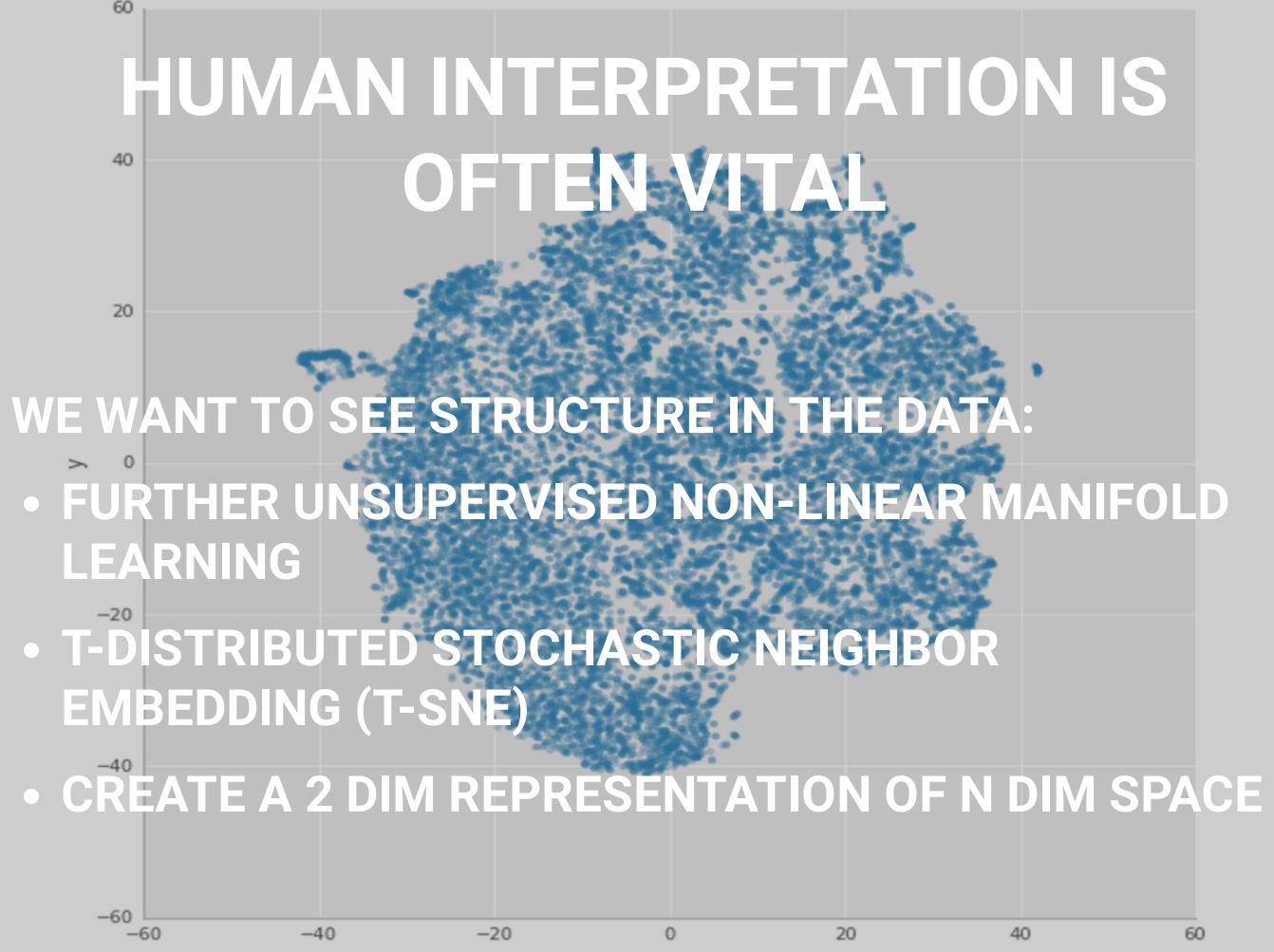
SINGULAR VALUE DECOMPOSITION

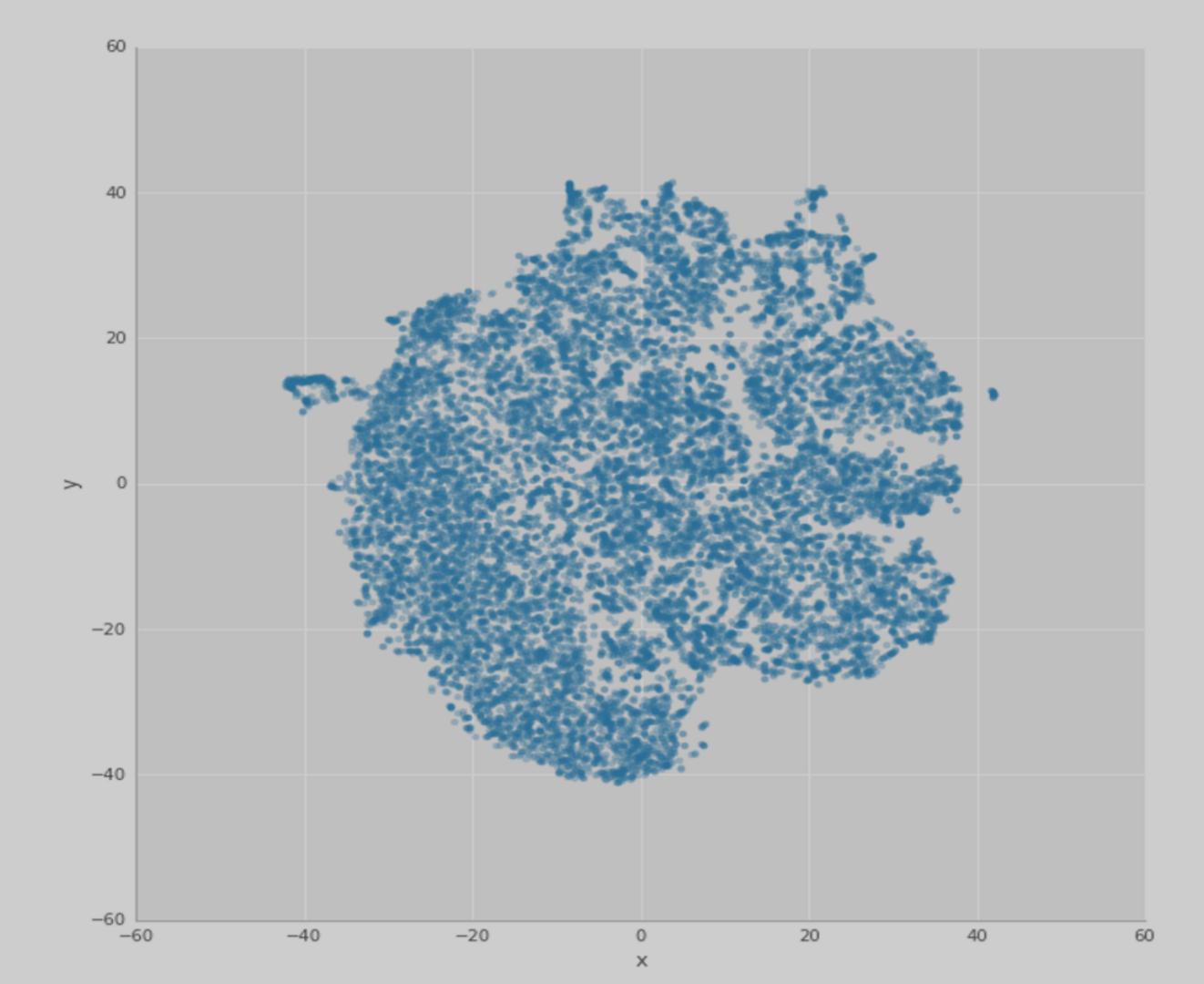
REDUCES THE SIZE OF THE PROBLEM BY DESCRIBING THE DATA IN A NEW SET OF AXES

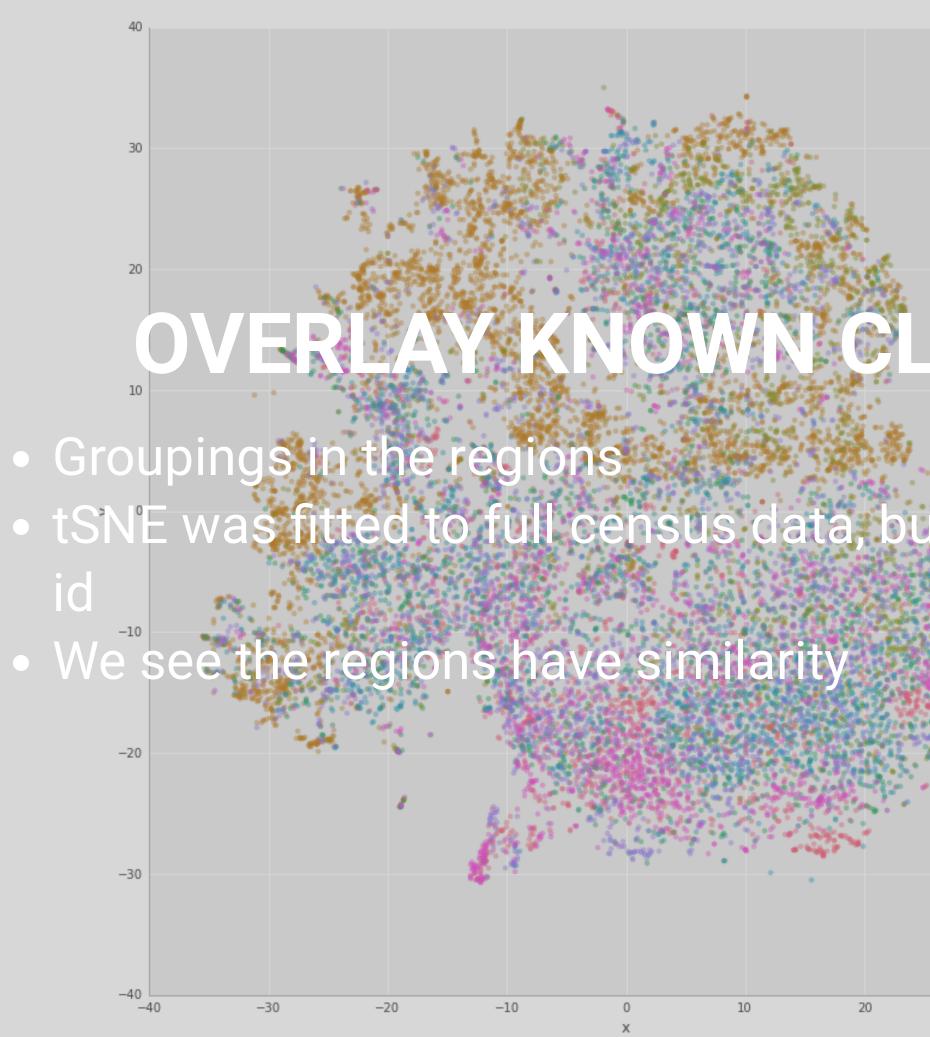
- FEATURES OFTEN PARTIALLY CORRELATE
- IF YOU KNOW ONE COLUMN, YOU PARTIALLY KNOW THE OTHER
- SO WE DESCRIBE BOTH USING A SINGLE COLUMN (WITH SOME MINOR LOSS OF INFO)







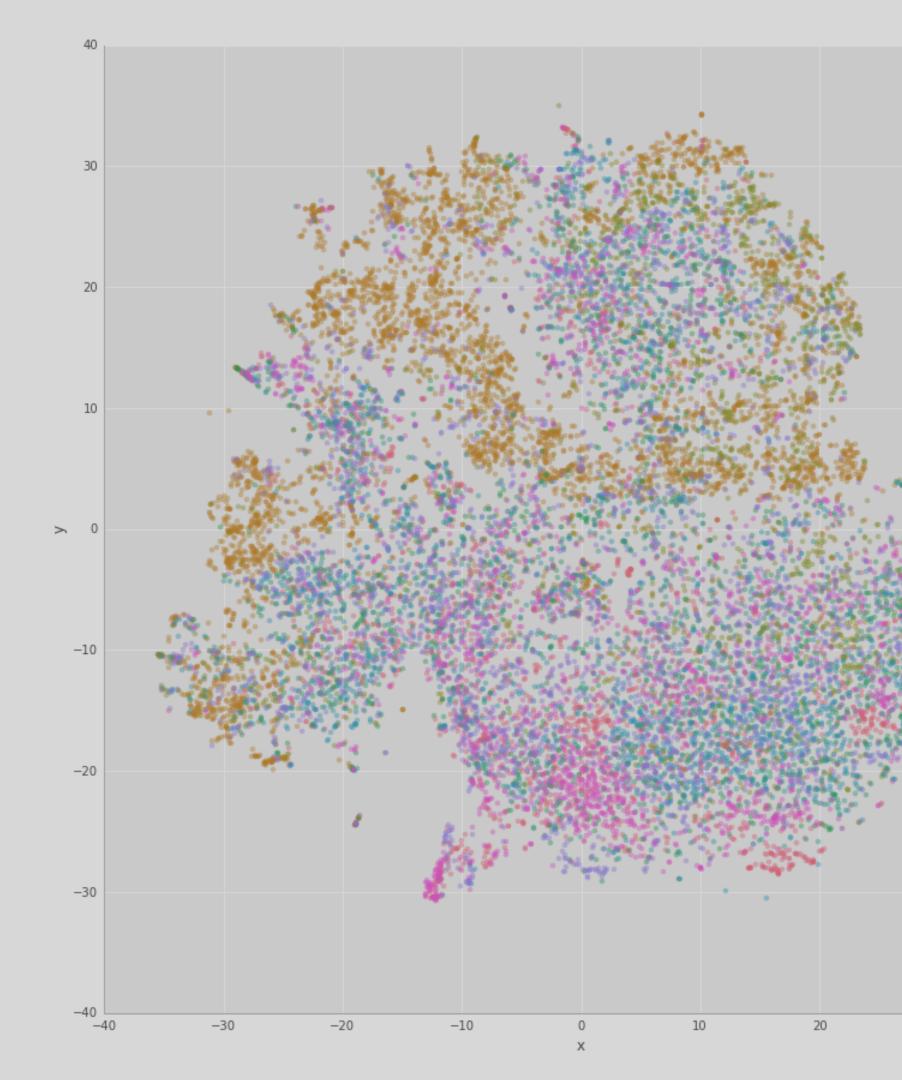




NCLASSES

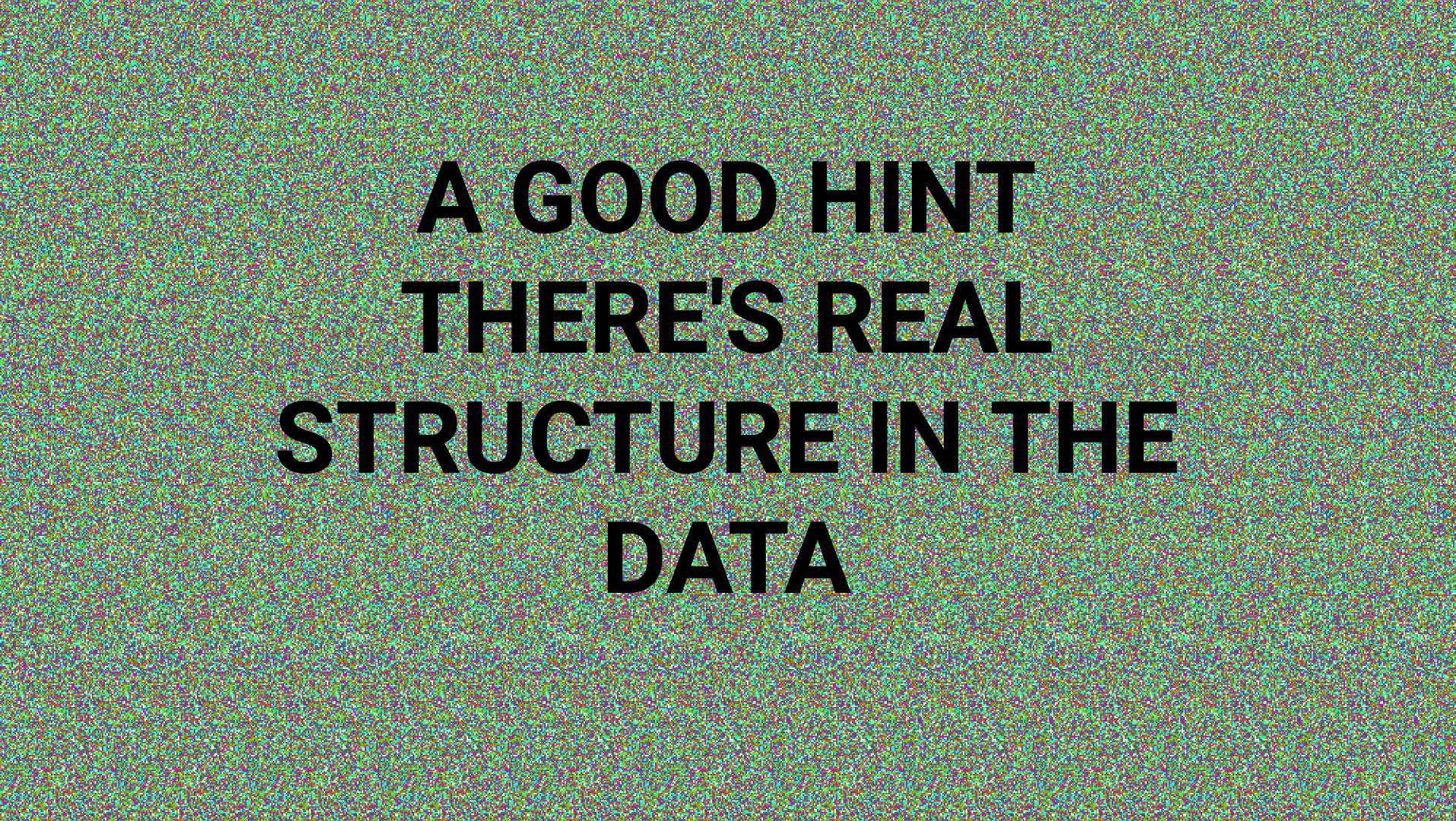
nuts3name

- Border
- Dublin
- Mid-East
 - Mi 🖽 it 🛛 🔹 Midland
- Midland
- South-East (IE) South-West (IE)
- South-wes
- West



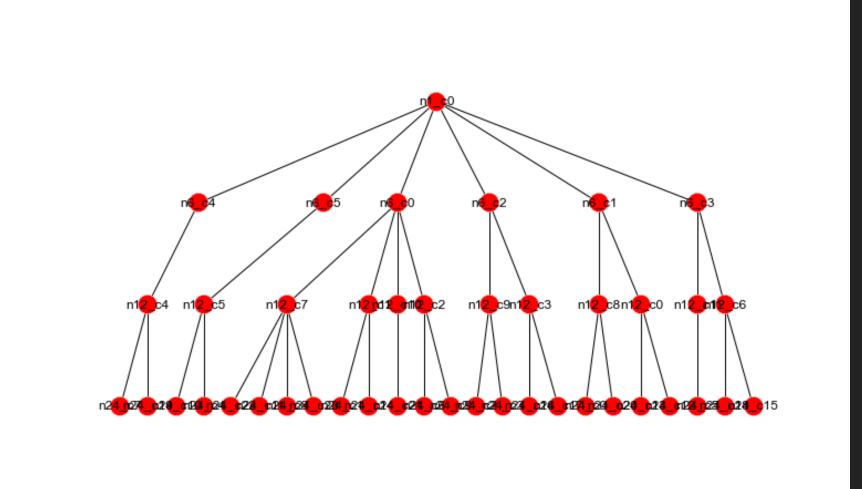
nuts3name

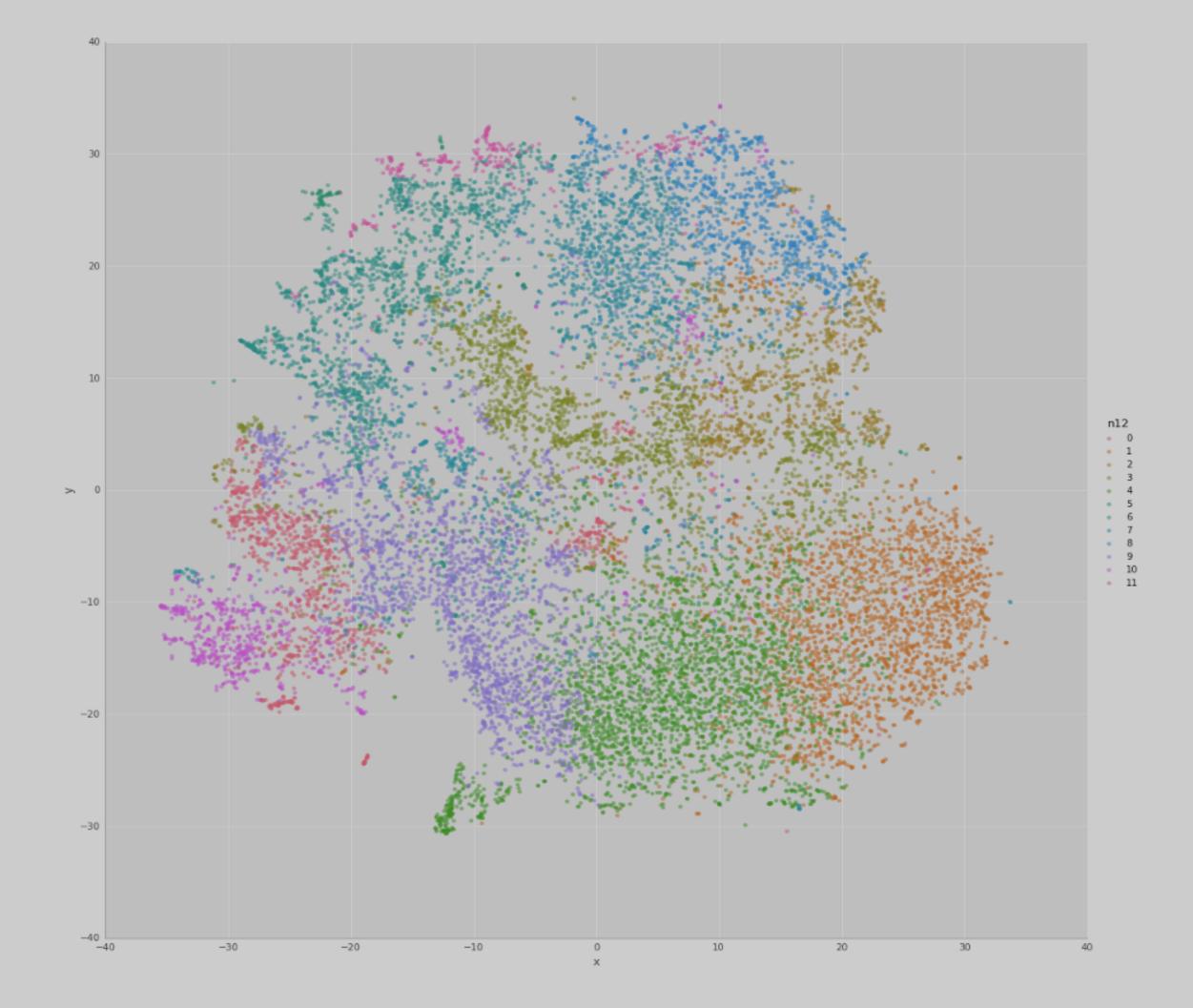
- Border
- Dublin
- Mid-East
- Mid-West
- Midland
- South-East (IE)
- South-West (IE)
- West



AGGLOMERATIVE HIERARCHICAL CLUSTERING

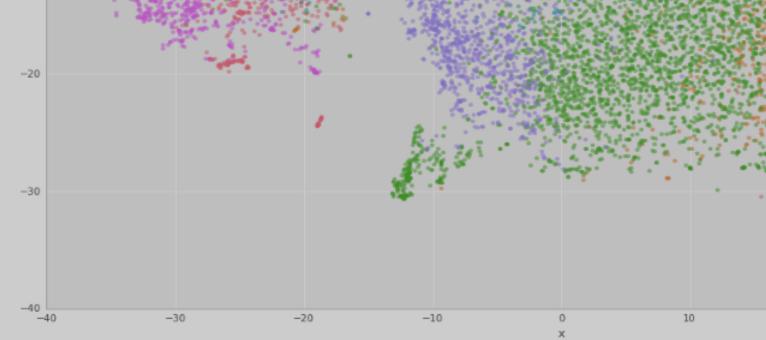
- Group nearby data points into progressively
 - larger clusters
- Get a nested hierarchy of clusters
- Choose your level





INTERESTING STRUCTURE!

- Clustering was entirely unsupervised
- i.e. determined only by the data itself
 Now we need to understand what the clusters mean





FIRST TRICK: MAP THEM

COLOR THE SMALL AREAS BY CLUSTER



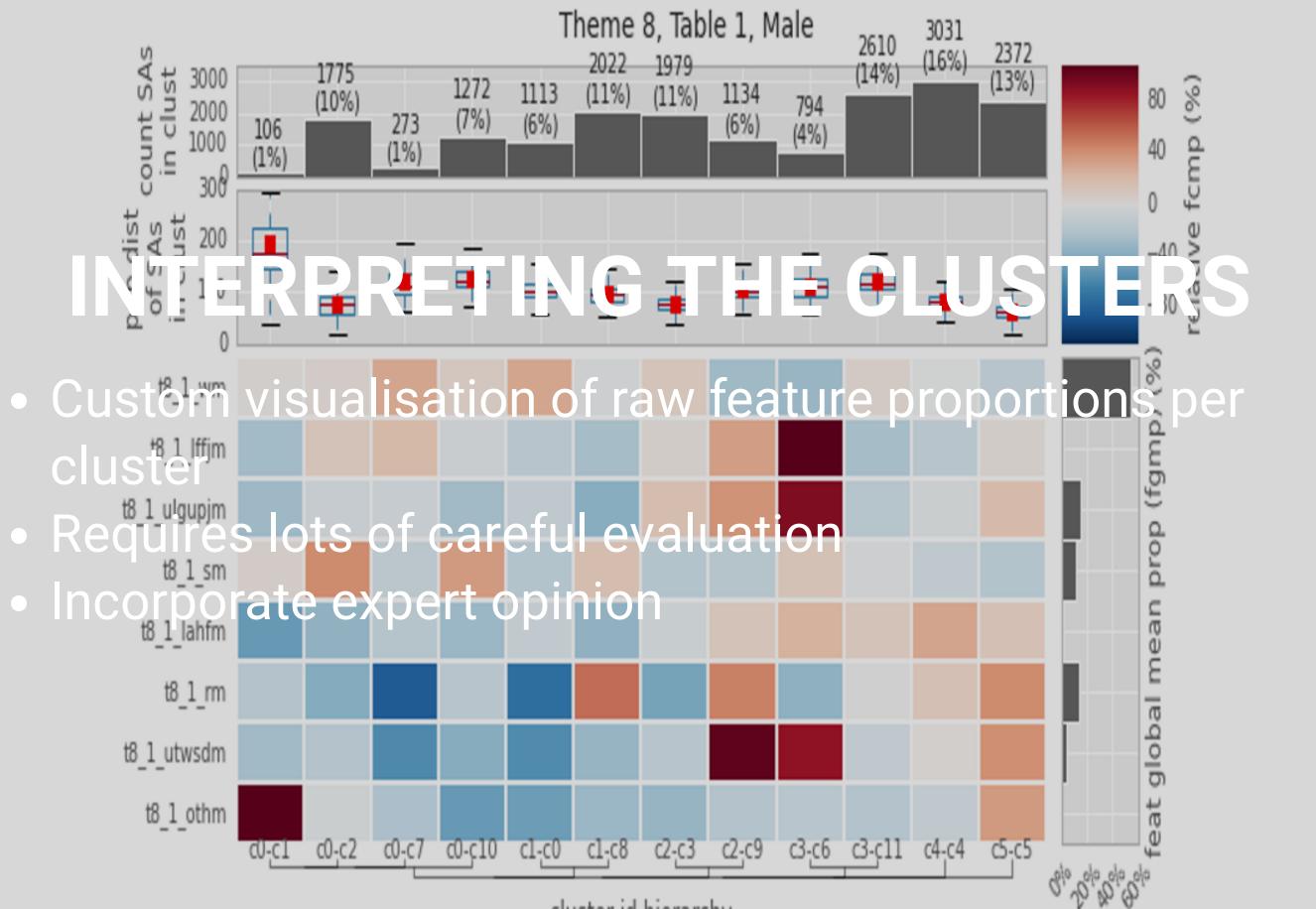


IF YOU LIVE IN IRELAND YOU CAN MAKE A GOOD GUESS...

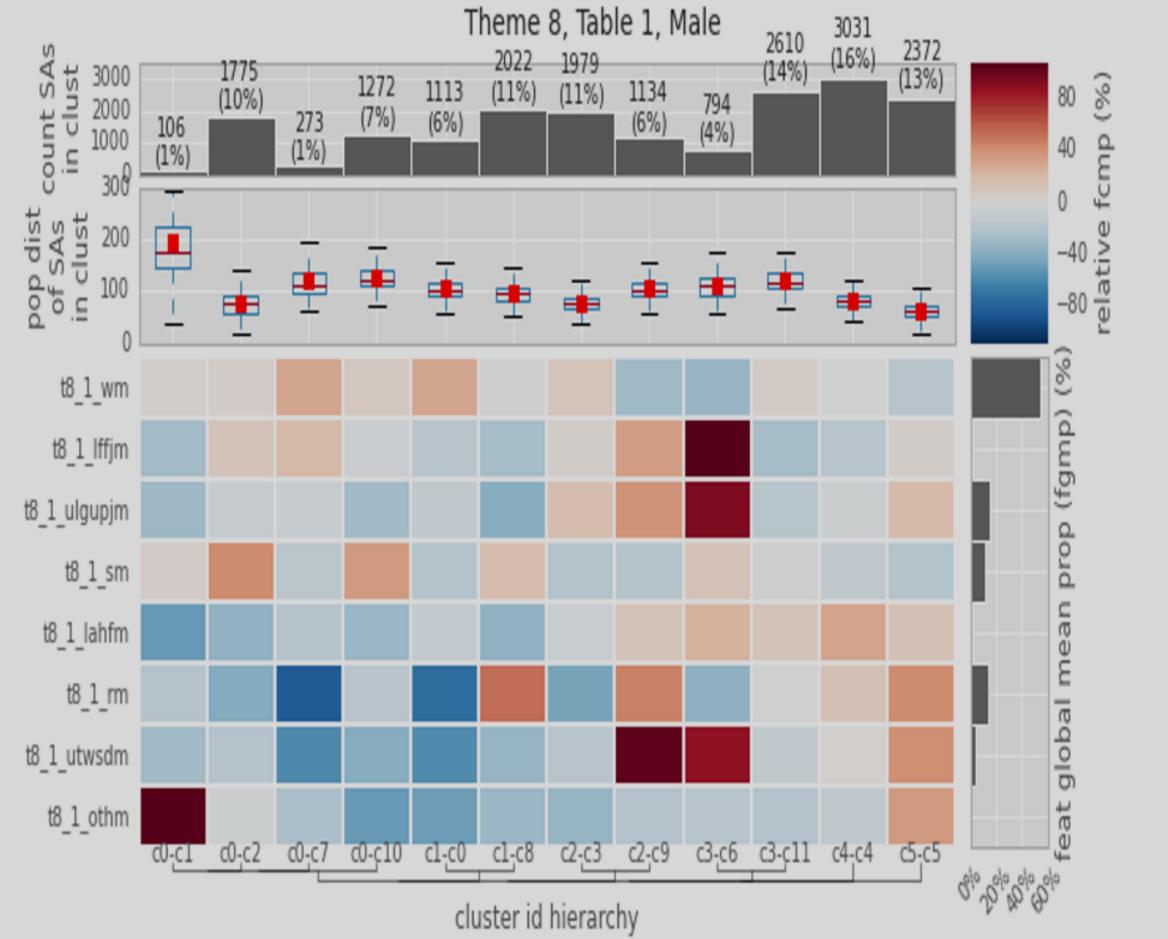
GUSTERS MEAN?

BUIGHERNEDO

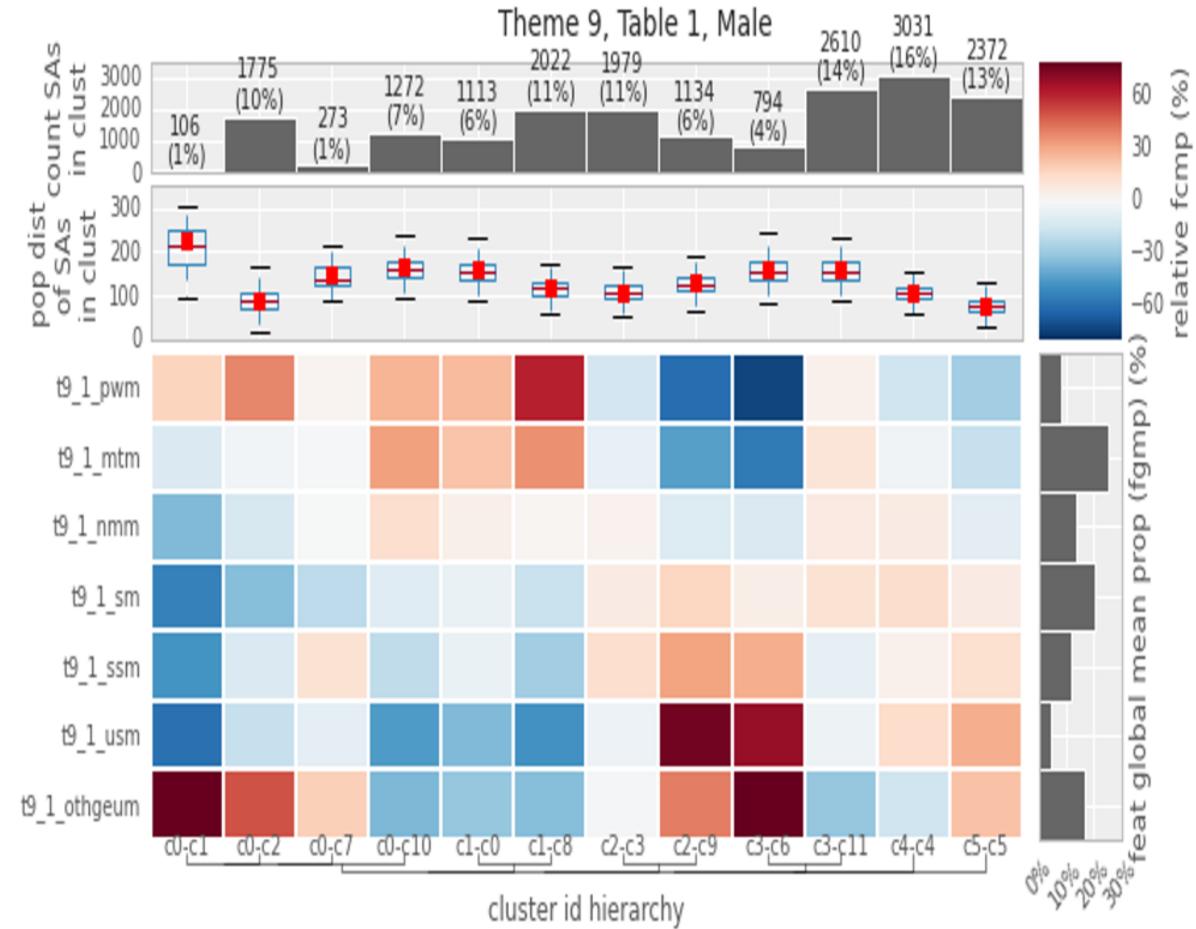
۲



cluster id hierarchy



cluster id hierarchy



cluster id hierarchy



INSUMARY DATA SCIENCE IS APPLICABLE por THROUGHOUT (THE INSURANCE -) n 1a 1 viding fir BUSINES protection agains death, loss, or damage. 1b the state called: insurance policy. the policy pecuniary amount of such protect tarra for anoth protections. 4. (an more

BUT MOST IMPORTANTLY REMEMBER IT'S ONLY A MODEL

QUESTIONS?

Nah, we're running outta time. Seriously.

THANKS FOR LISTENING