

Doing Visual Data Science

Foundations, Techniques and Practice

Cagatay Turkay
Professor,
Centre for Interdisciplinary Methodologies
University of Warwick

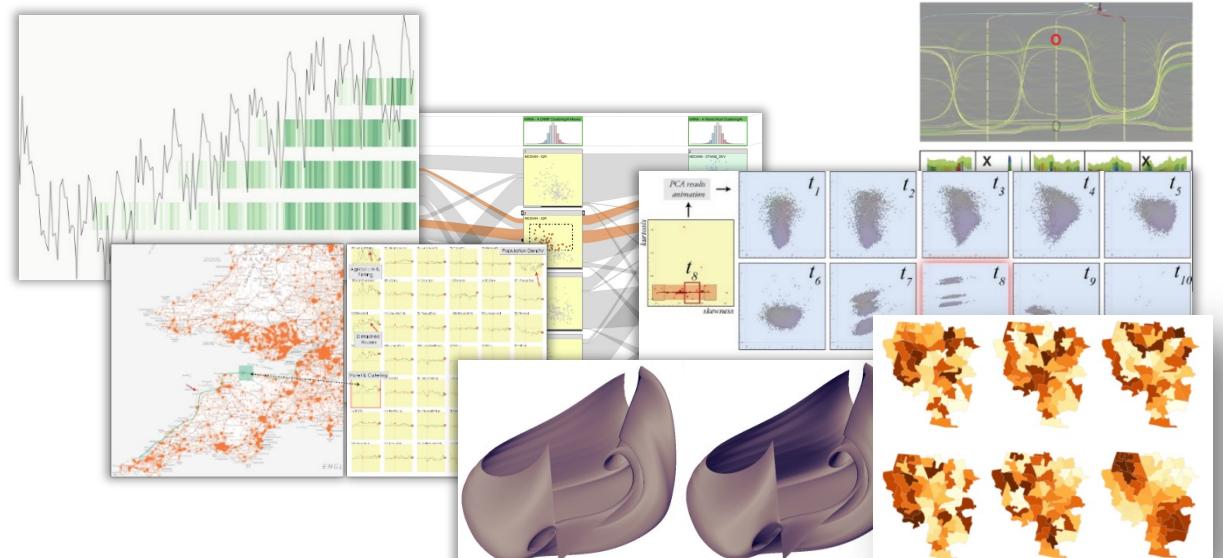
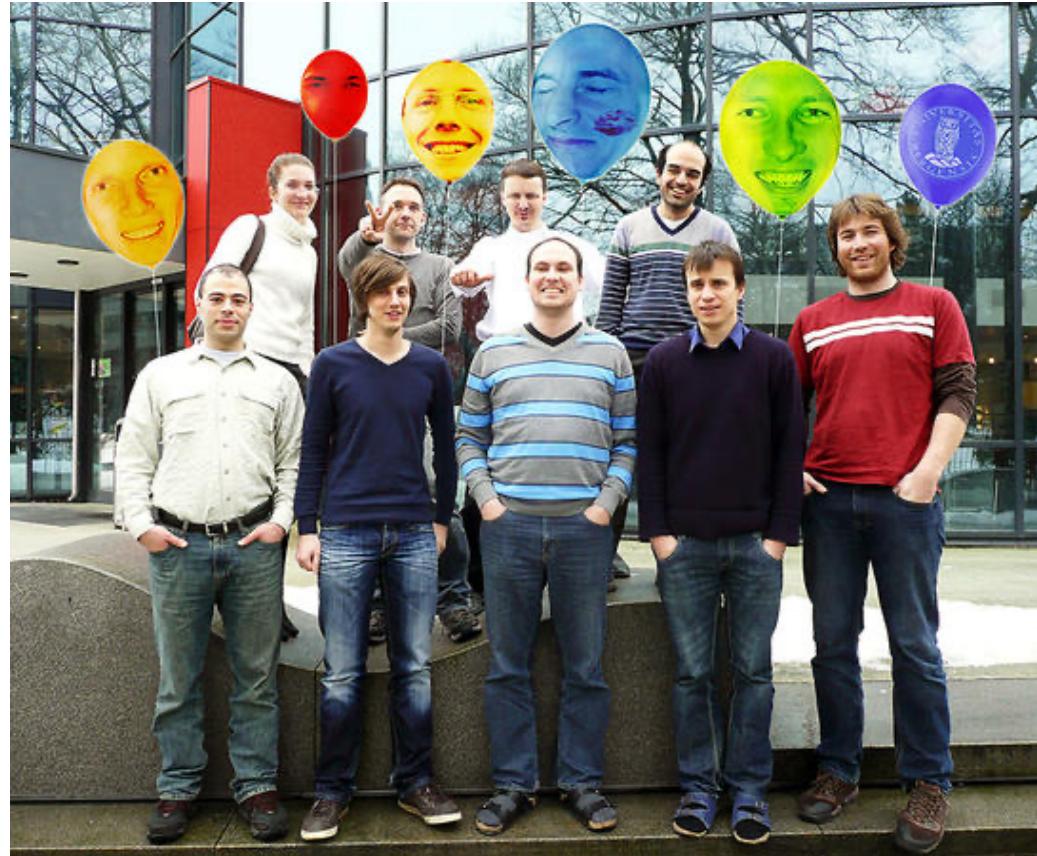


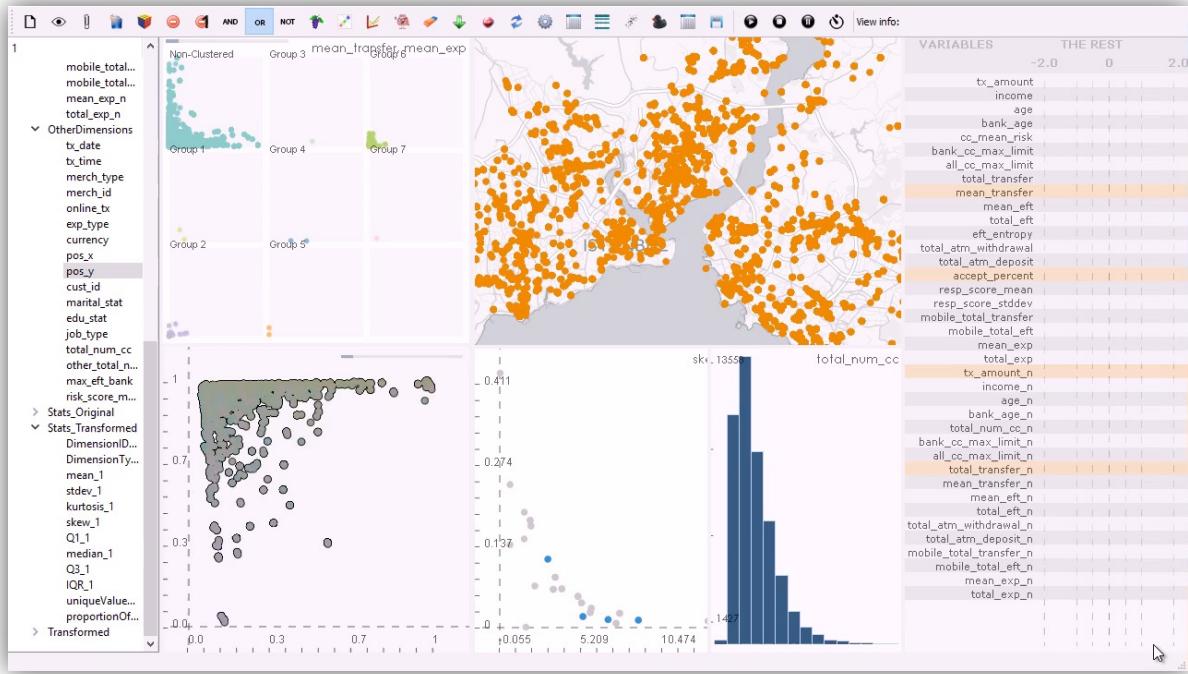
Professor @ CIM at University of Warwick

Senior Lecturer in at giCentre, City, University of London (2014 - 2019)

PhD in Visualisation, University of Bergen, NO

MSc & BSc in Computer Science (Sabanci University, TR & METU, TR)



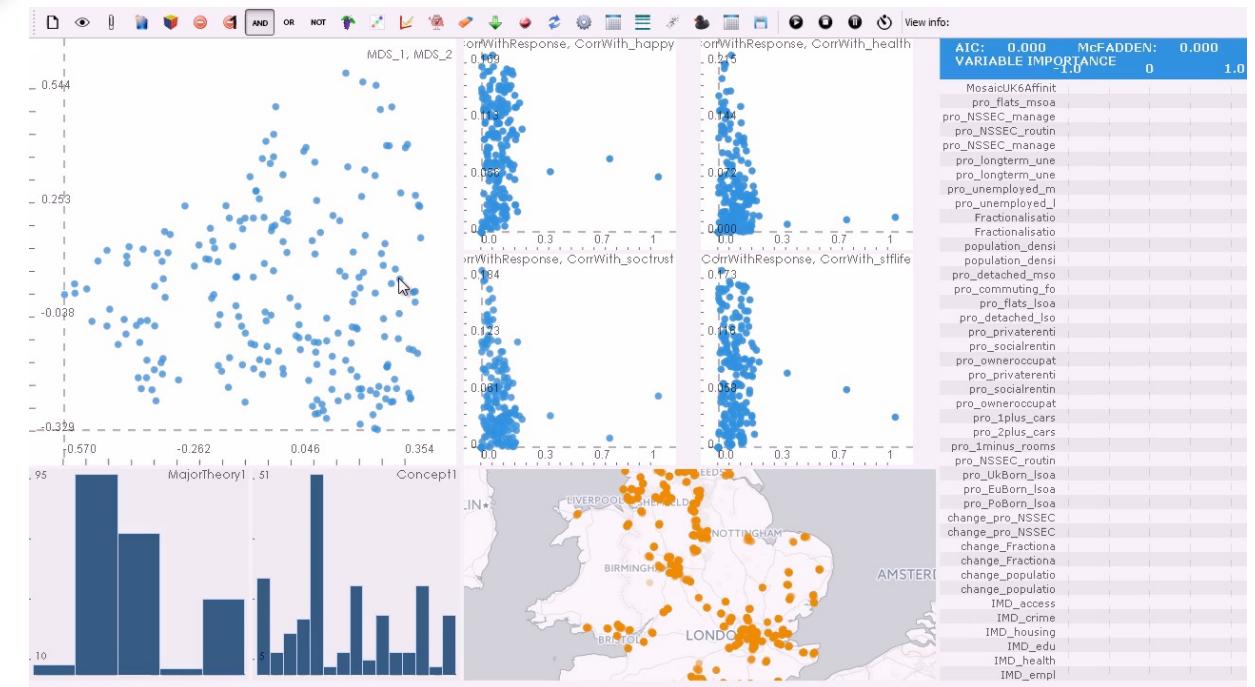


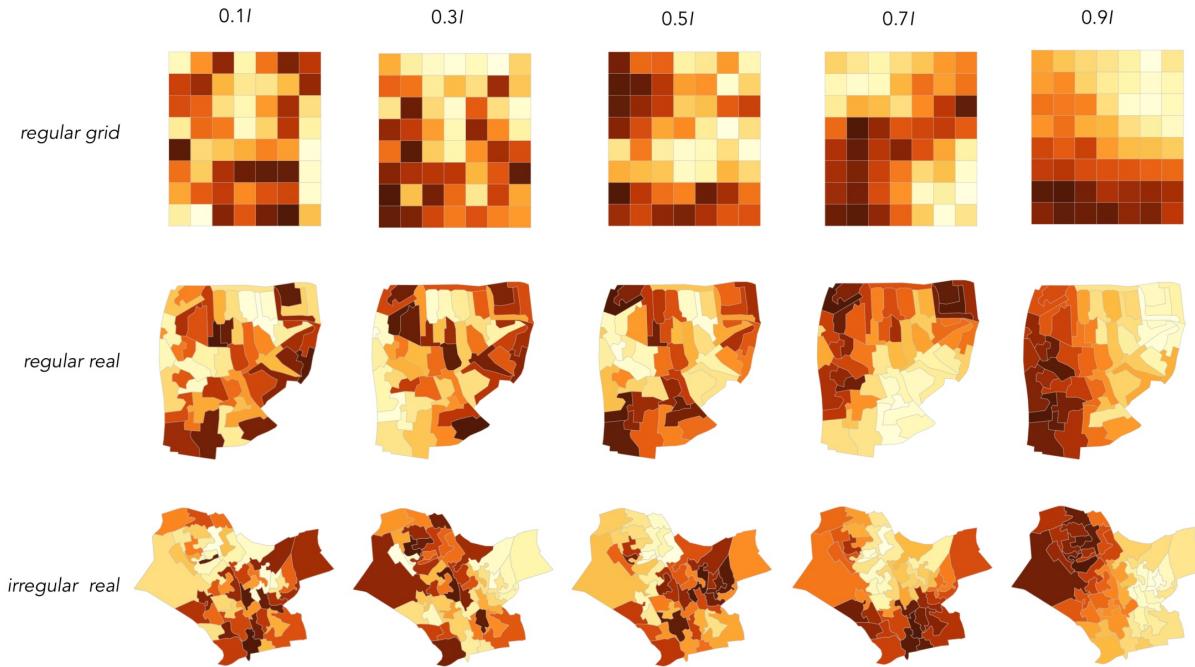
Enhancing a Social Science Model-building Workflow with Interactive Visualisation

Turkay, C., Slingsby, A., Lahtinen, K., Butt, S., & Dykes, J., ESANN 2016 (& Neurocomputing 2017)

Designing Progressive and Interactive Analytics Processes for High-Dimensional Data Analysis

Turkay, C., Kaya, E., Balcisoy, S., Hauser, H., IEEE TVCG 2017



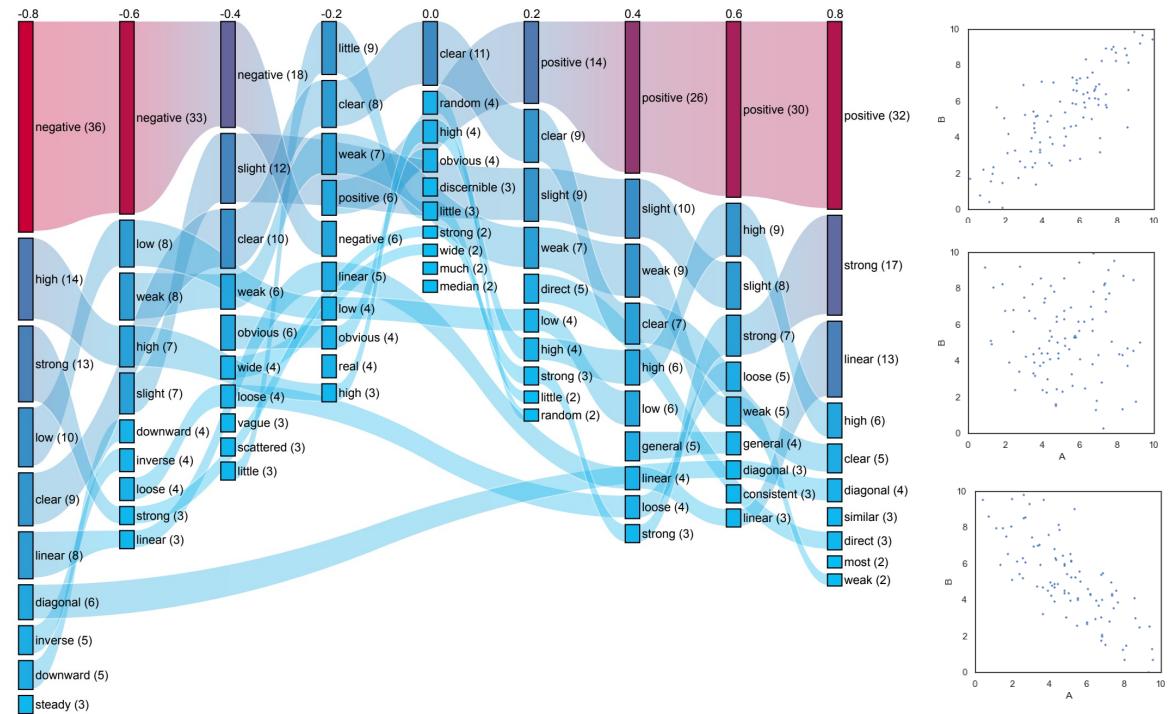


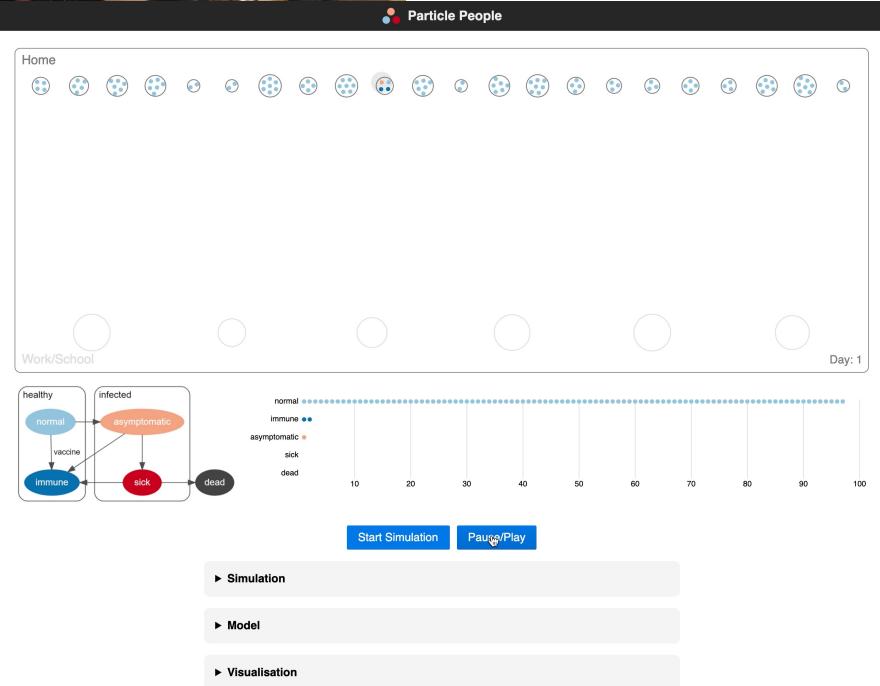
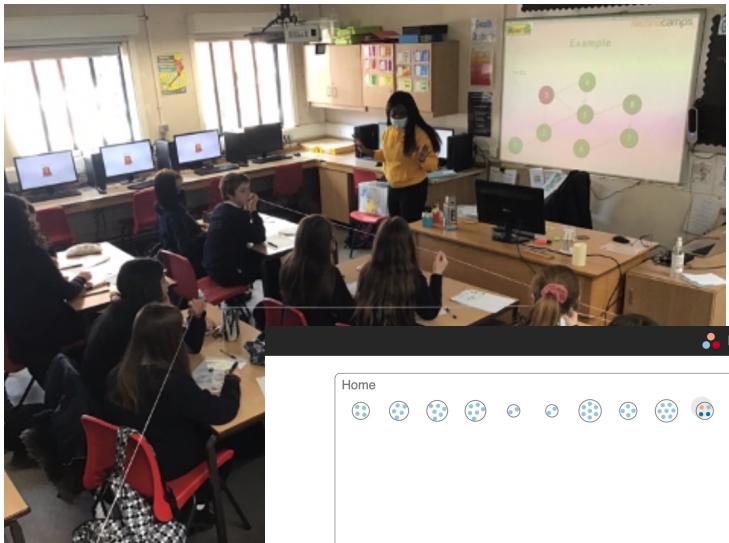
Map LineUps: effects of spatial structure on graphical inference

Beecham, R., Dykes, J., Meulemans, W., Slingsby, A., Turkay, C. & Wood, J. (2016).
IEEE Transactions on Visualization and Computer Graphics.

Words of Estimative Correlation: Studying Verbalizations of Scatterplots

Henkin, R., Turkay, C.(2021).
IEEE Transactions on Visualization and Computer Graphics.





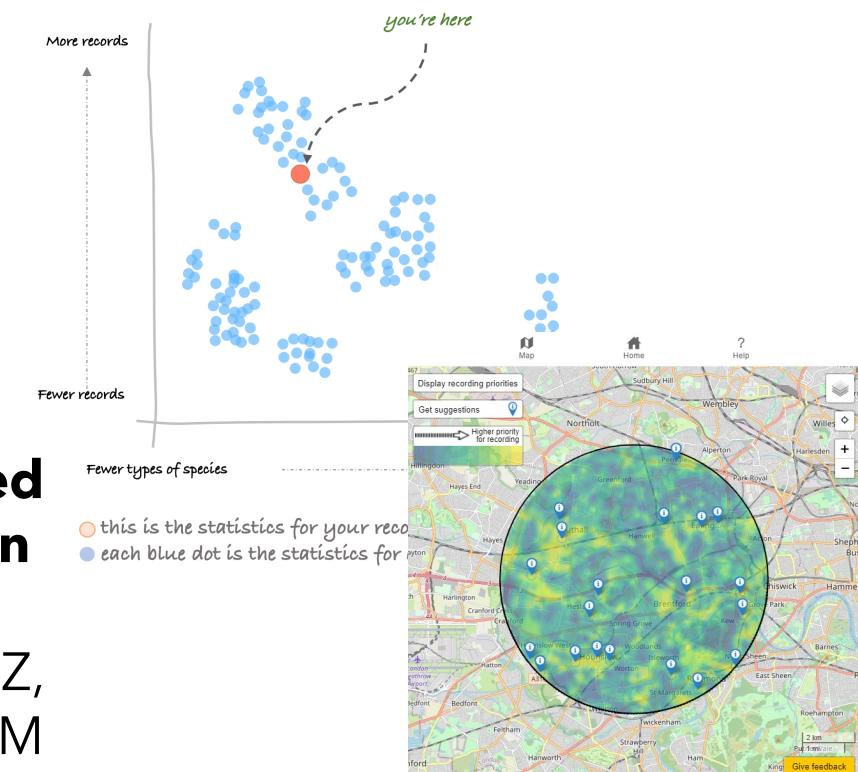
'Data stories': inspiring citizen scientists with personalised automated data-driven narrative visualisation

Rolph S, Turkay C, August T, Pateman R, Baird K, Randle Z,
Siddharthan A, Bowler A, Hutchinson R and Pocock M
(in preparation)

From Asymptomatics to Zombies: Visualization-Based Education of Disease Modeling for Children.

McNeill, G., Sondag, M., Powell, S., Asplin, P., Turkay, C., Moller, F. and Archambault, D, 2023 CHI Conference on Human Factors in Computing Systems

See how diverse recorders you and the DECIDE community are ...



VDS – Our itinerary for today

- Visual data science?
- Why visual?
- A bit of historical context
- VDS practice & techniques

- On to DIY...

HOW DO YOU UNDERSTAND “DATA SCIENCE”?

Is it a discipline, a bunch of techniques, a new name for statistics, ...?

.. *data science as a **process***

... starting with a problem/question that can be **approached with data**, and involves **iteratively collecting, cleaning, analysing** and developing models in dialogue with theories and assumptions about the phenomena being investigated while communicating along this iteration

.. visual data science as a(n)

method, approach & way of thinking

... that is underpinned by visual and interactive computing methods ...

... to interrogate, reason about and model a phenomena of interest through data and computational artefacts ..

.. through these 2 days, demonstrate & experience ...

VISUALISATION
TO “**BETTER**”
DATA SCIENCE

RICH,
ENGAGED,
REFLECTIVE,
INFORMED

*dialogue with
data & models*

**Why we need a
rich, engaged, reflective and
informed
dialogue with data & models?**

Data & Models will always tell you a story ... no matter what

Storks Deliver Babies ($p = 0.008$)

KEYWORDS:

Teaching;
Correlation;
Significance;
 p -values.

Robert Matthews

Aston University, Birmingham, England.
e-mail: rajm@compuserve.com

Summary

This article shows that a highly statistically significant correlation exists between stork populations and human birth rates across Europe. While storks may not deliver babies, unthinking interpretation of correlation and p -values can certainly deliver unreliable conclusions.

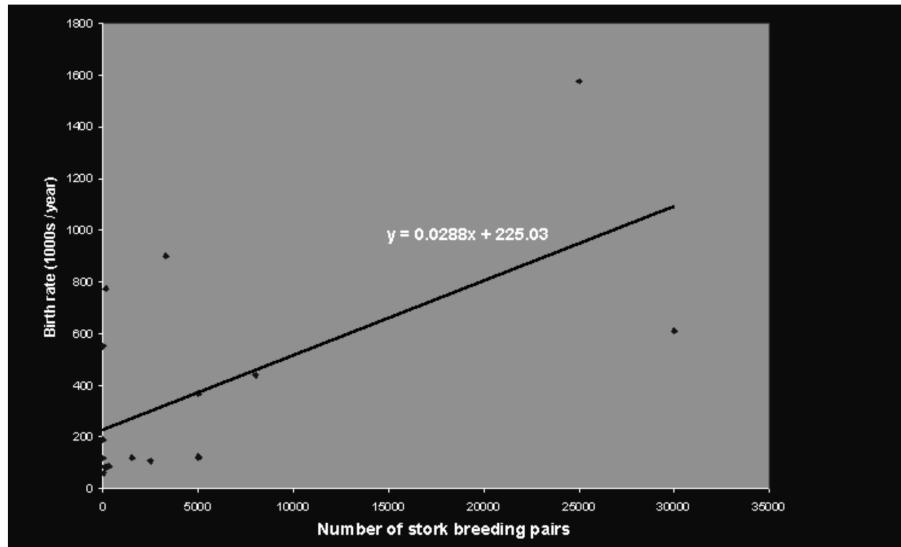


Fig 1. How the number of human births varies with stork populations in 17 European countries.

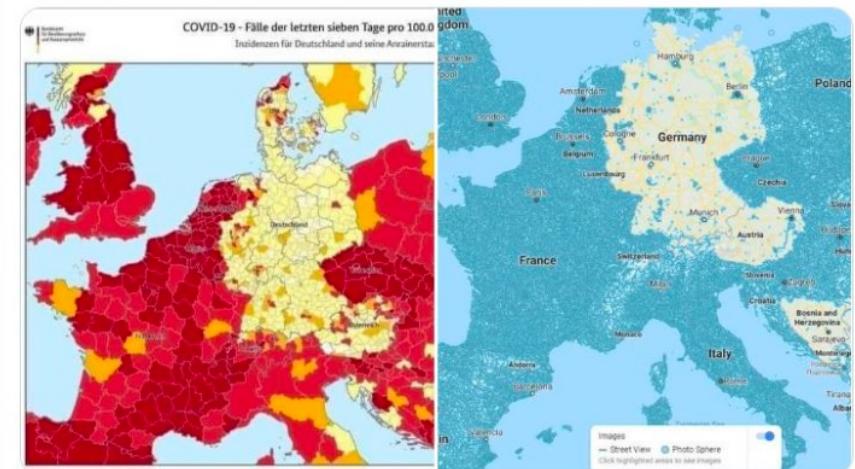
← Tweet



Luuk van der Meer
@LuukvanderMeer

The evidence cannot be clearer: Google Street View causing COVID-19. #correlation #causation #conspiracy #COVID19

Przetłumacz Tweeta



7:31 AM · 17 paź 2020 · Twitter for Android

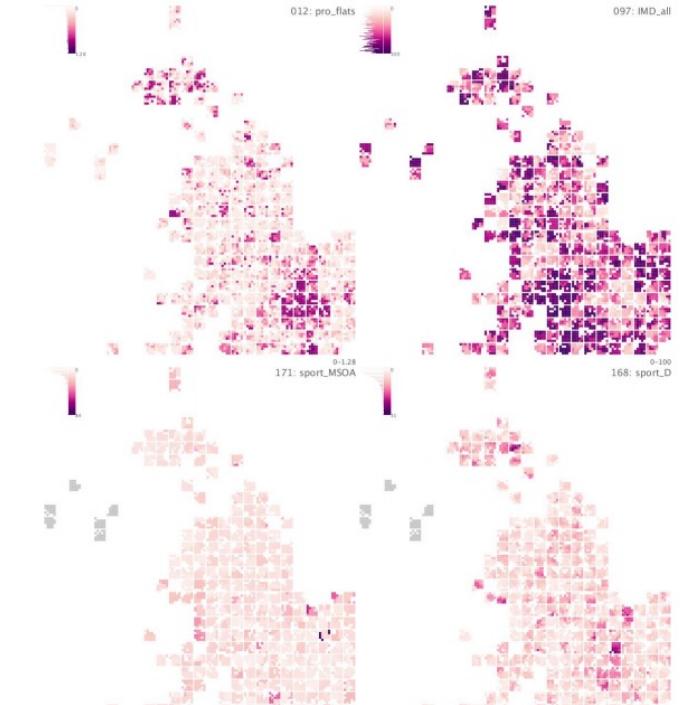
1,4 tys. Tweetów podanych dalej 113 Cytatów z Tweetów 5,5 tys. Polubień

<https://twitter.com/LuukvanderMeer/status/1317352613592137728>

Phenomena is complex .. need to understand data & models ..

*“We (social scientists) need (data-based) models that we can **understand and explain** so that we can defend them to our peers in **full confidence.**”*

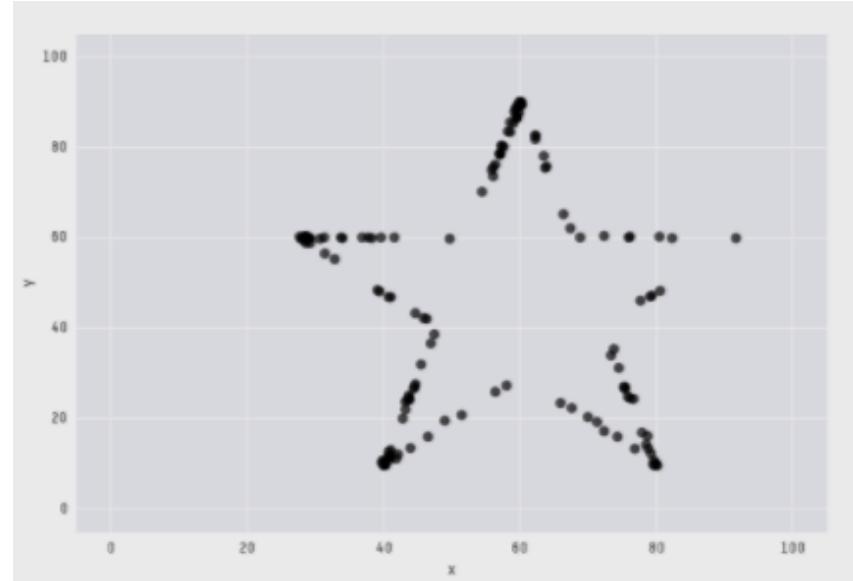
A quote from collaborators at our AddResponse project on data-based models

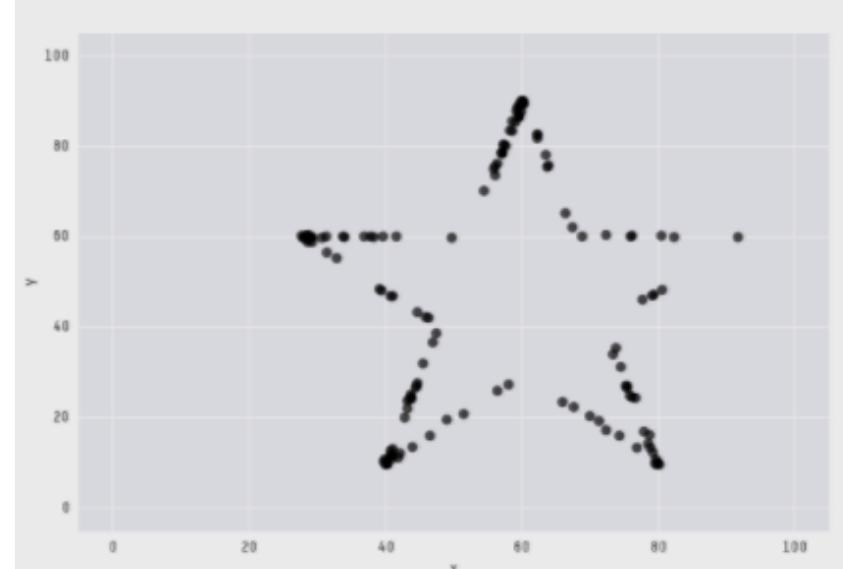


From: Lahtinen, K. et al. (2015). Informing Non-Response Bias Model Creation in Social Surveys with Visualisation. Poster VIS 2015

WHY VISUALISE?

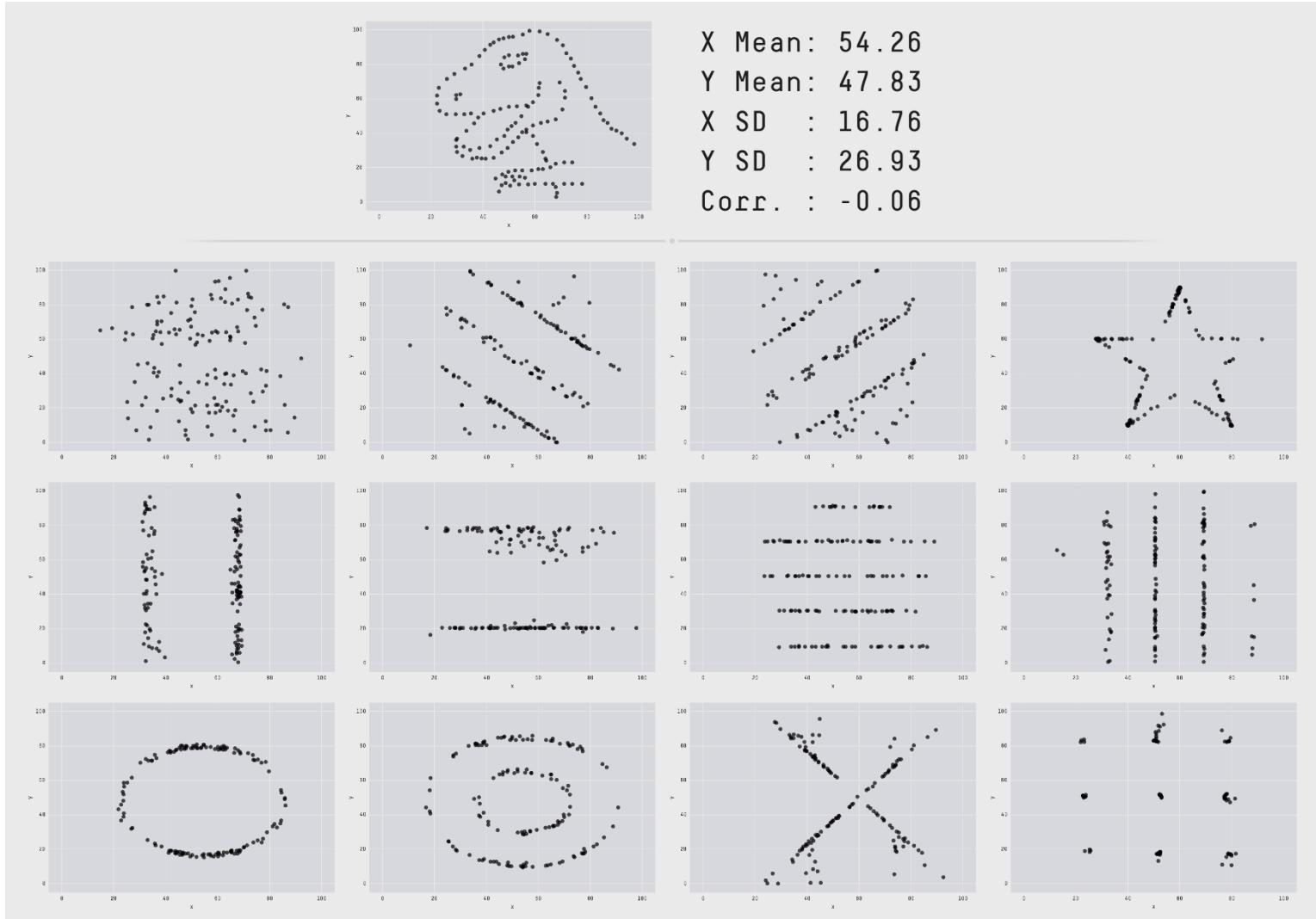
Do you see any similarities?





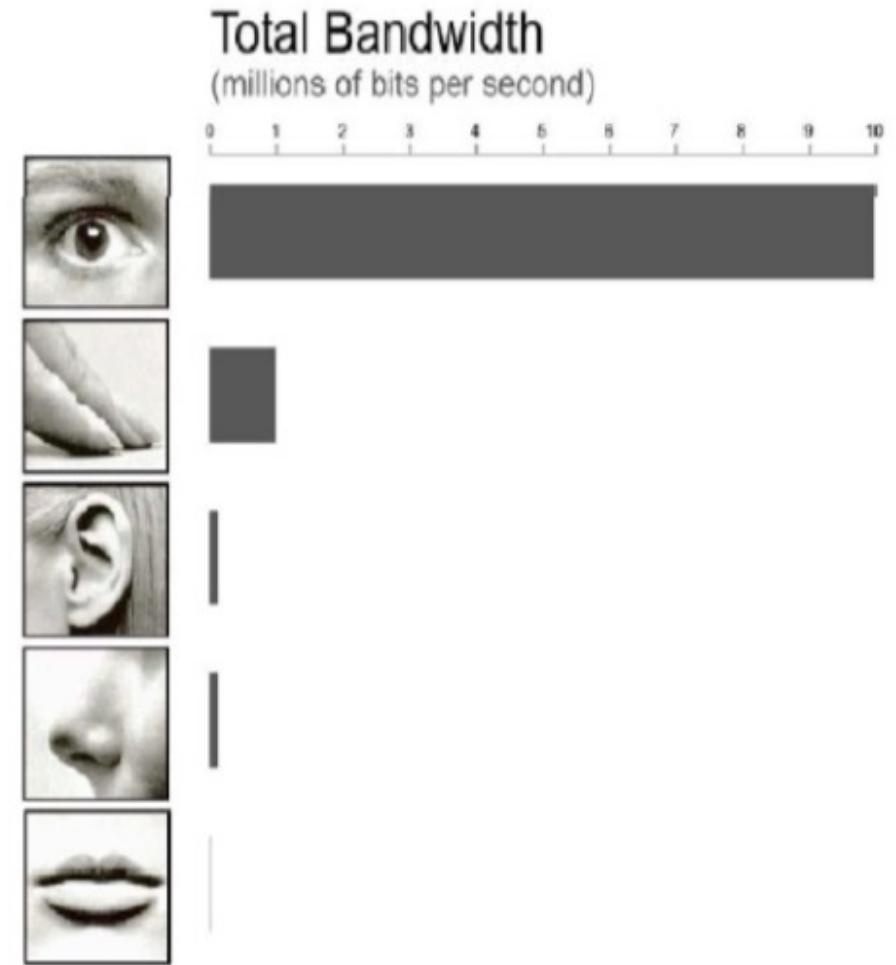
X Mean: 54.26
Y Mean: 47.83
X SD : 16.76
Y SD : 26.93
Corr. : -0.06

The Datasaurus Dozen



Why Visual?

- Figures are **richer**; provide more information with less clutter and in less space.
- Figures provide the *gestalt* effect: they give an overview; **make structure more visible**.
- Figures are more **accessible**, easier to understand, **faster to grasp**, more **comprehensible**, more **memorable**, more **fun**, and **less formal**.



List & slide adapted from: [Stasko et al. 1998 & Alexander Lex]

E N C O D I N G

D E C O D I N G

A C T I O N

$x_1, y_1,$
 $x_2, y_2 \dots$



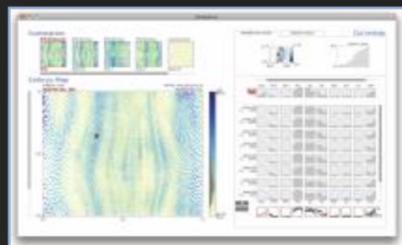
Engage



“Wow, X & Y
looks
amazing!”

“I need to
find out
more!”

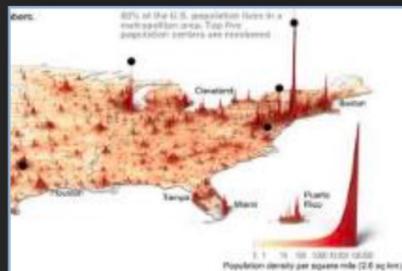
Explore



“I wonder how
x relates to y”

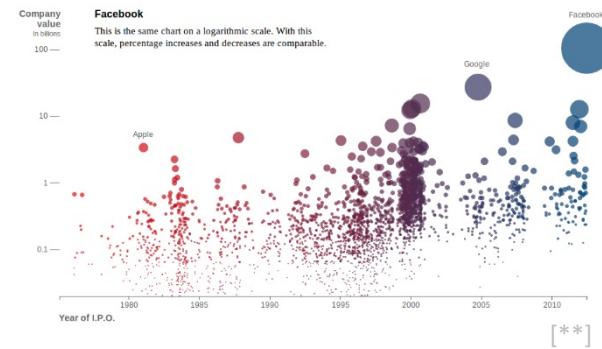
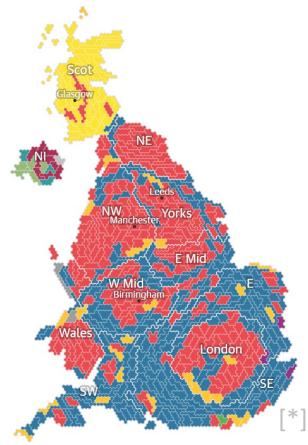
“Maybe z is
important?”

Explain



“X does y”

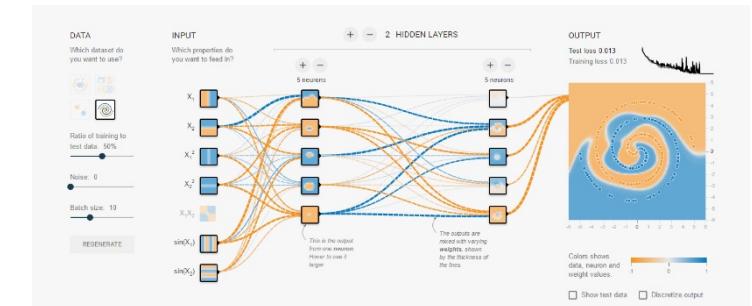
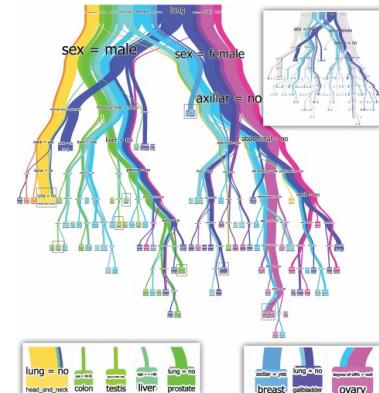
“if I do x
then...”



VIS for communication



VIS for analysis



<http://playground.tensorflow.org>

[Elzen & Wijk, 2010]

Sources:

[*] <http://www.theguardian.com/politics/ng-interactive/2015/apr/20/election-2015-constituency-map>

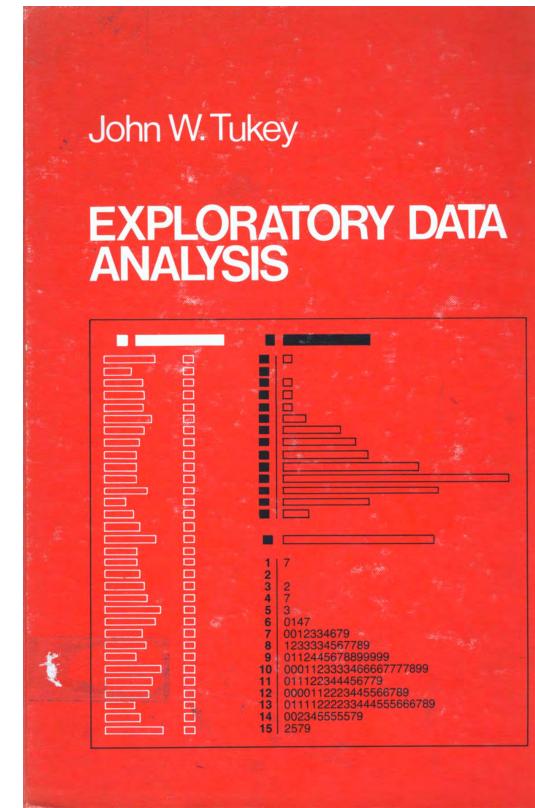
[**] <http://www.nytimes.com/interactive/2012/05/17/business/dealbook/how-the-facebook-offering-compares.html>

A BIT OF HISTORICAL CONTEXT?



Exploratory data analysis (EDA)

is a data-driven analysis approach where properties, structures, and patterns in the data are iteratively interpreted, leading to new understandings on the phenomena being investigated and avenues for further analytical enquiry.



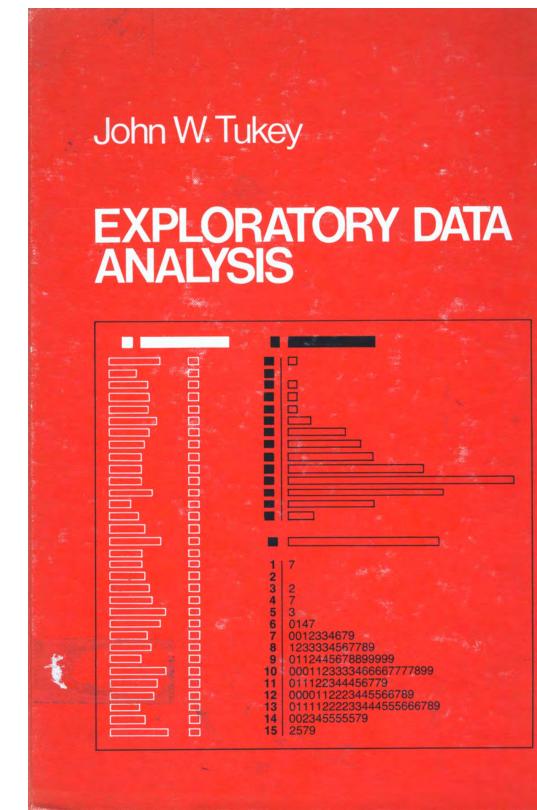
Exploratory data analysis (EDA)

An approach/tradition to data analysis where there is an:

- "(a) an **emphasis on** the substantive **understanding of data** that address the broad question of "what is going on here?"*
- (b) an **emphasis on graphic representations** of data;*
- (c) a focus **on tentative model building and hypothesis generation** in an **iterative process** of model specification, residual analysis, and model respecification;*
- (d) use of robust measures, re-expression, and subset analysis;*
- (e) positions of **skepticism, flexibility, and ecumenism** regarding which methods to apply."*

From John Behrens, 1997

Behrens, J.T., 1997. Principles and procedures of exploratory data analysis. Psychological Methods, 2(2), p.131.



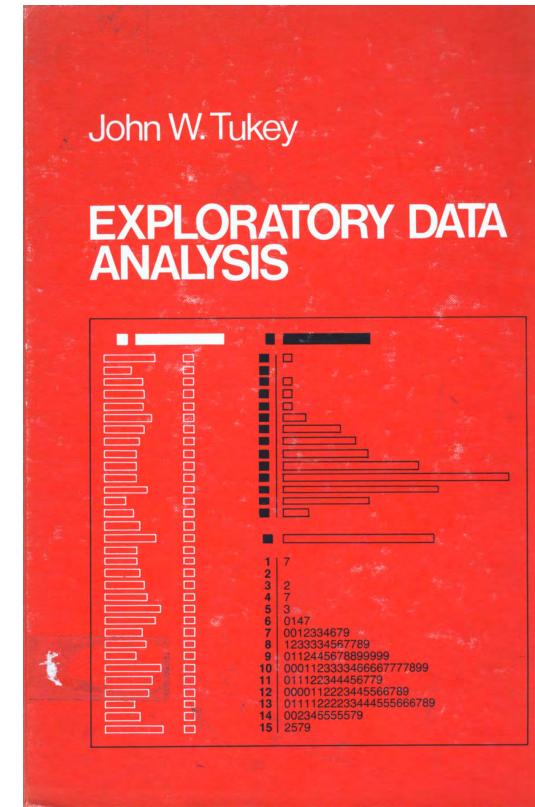
"Exploratory data analysis can never be the whole story, but nothing else can serve as the foundation stone--as the first step" (Tukey, 1977, p. 3).



"The simple graph has brought more information to the data analyst's mind than any other device."

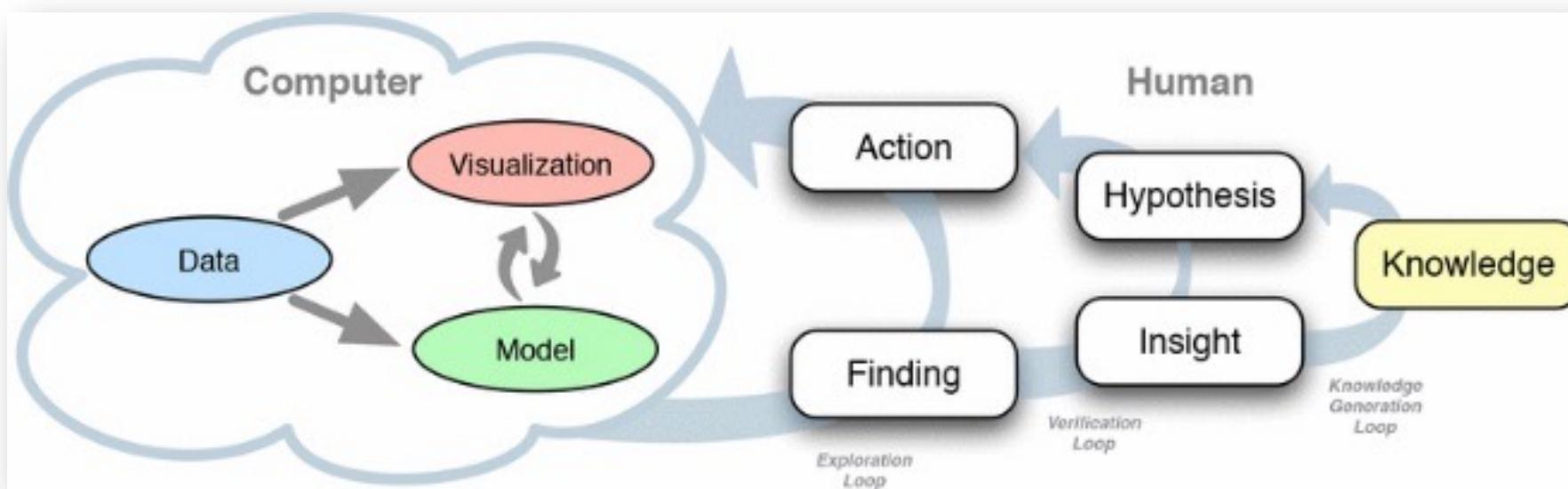
"Far better an approximate answer to the right question, which is often vague, than an exact answer to the wrong question, which can always be made precise."

THE FUTURE OF DATA ANALYSIS¹
BY JOHN W. TUKEY
Princeton University and Bell Telephone Laboratories

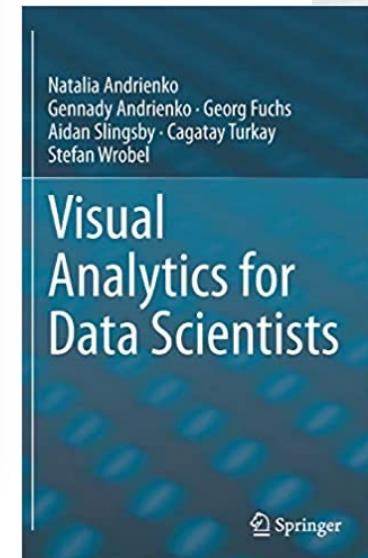
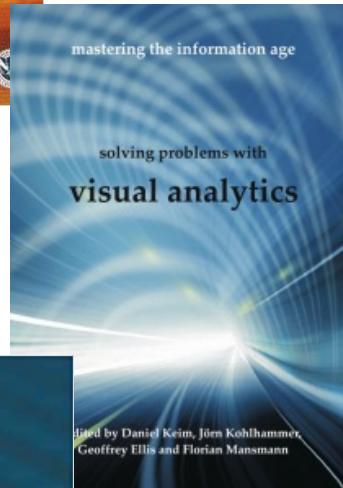
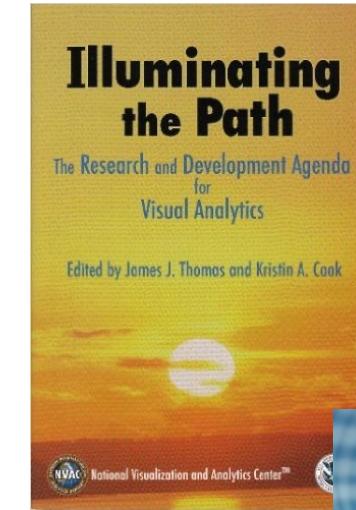


Visual Analytics

Combining visual and computational analysis:
human strengths & computing



[SACHA et al., 2014]



.. *visual data science as a(n)*

approach / discipline / way of thinking

... that is underpinned by visual and interactive computing methods ...

... to interrogate, reason about and model a phenomena of interest through data and computational artefacts ..

VDS: PRACTICE & TECHNIQUES

(a non-comprehensive tour!)

Four key practices in Visual Data Science

... facilitated by interaction and visualisation

P1: interrogate, relate, compare **variation, co-variation, deviation**

P2: **interact** with data and computational tools

P3: generate **new perspectives** to “see” data and models

P4: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

Four key practices in Visual Data Science

... facilitated by interaction and visualisation

P1: interrogate, relate, compare **variation, co-variation, deviation**

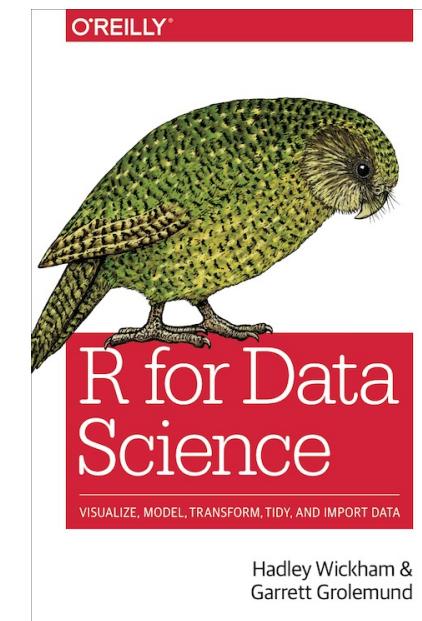
P2: **interact** with data and computational tools

P3: generate **new perspectives** to “see” data and models

P4: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

“... two types of questions will always be useful for making discoveries within your data. You can loosely word these questions as:

- 1. What type of **variation** occurs **within my variables**?*
- 2. What type of **covariation** occurs **between my variables**?”*



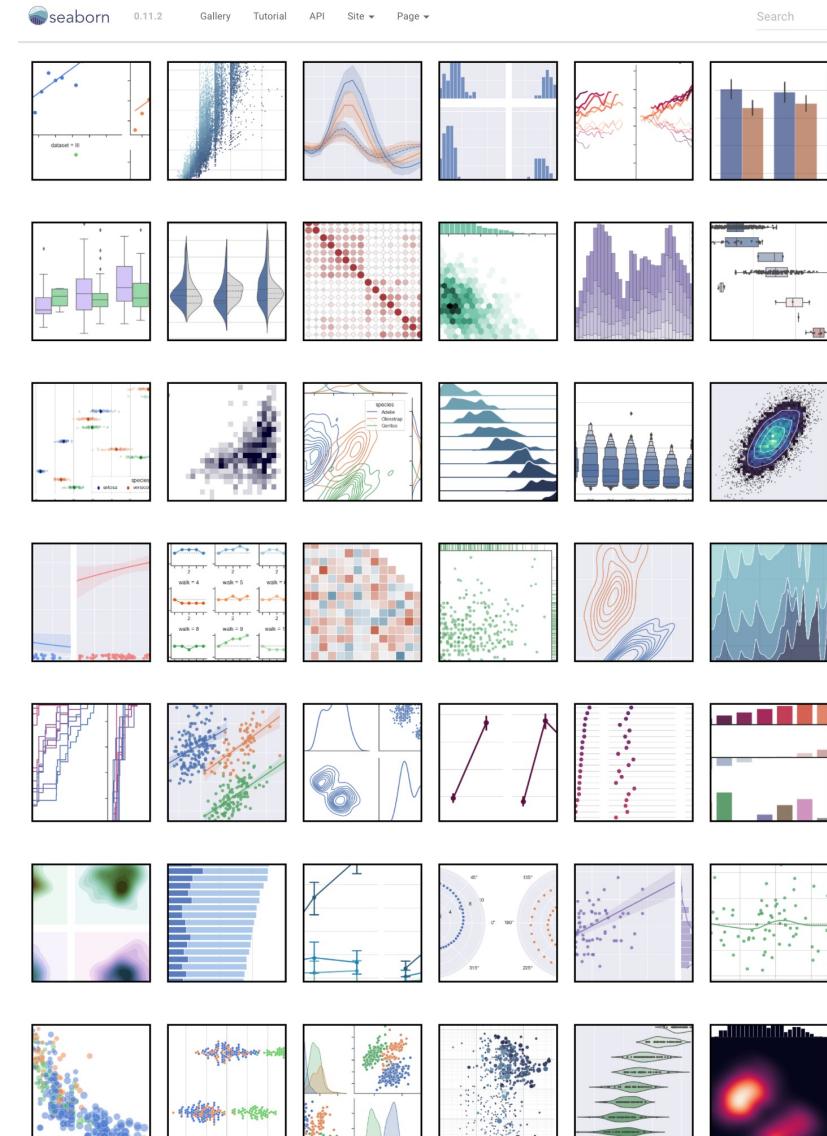
From: <https://r4ds.had.co.nz/exploratory-data-analysis.html>

Hadley Wickham &
Garrett Grolemund

Many forms of visualizations for exploring variation, co-variation, deviation



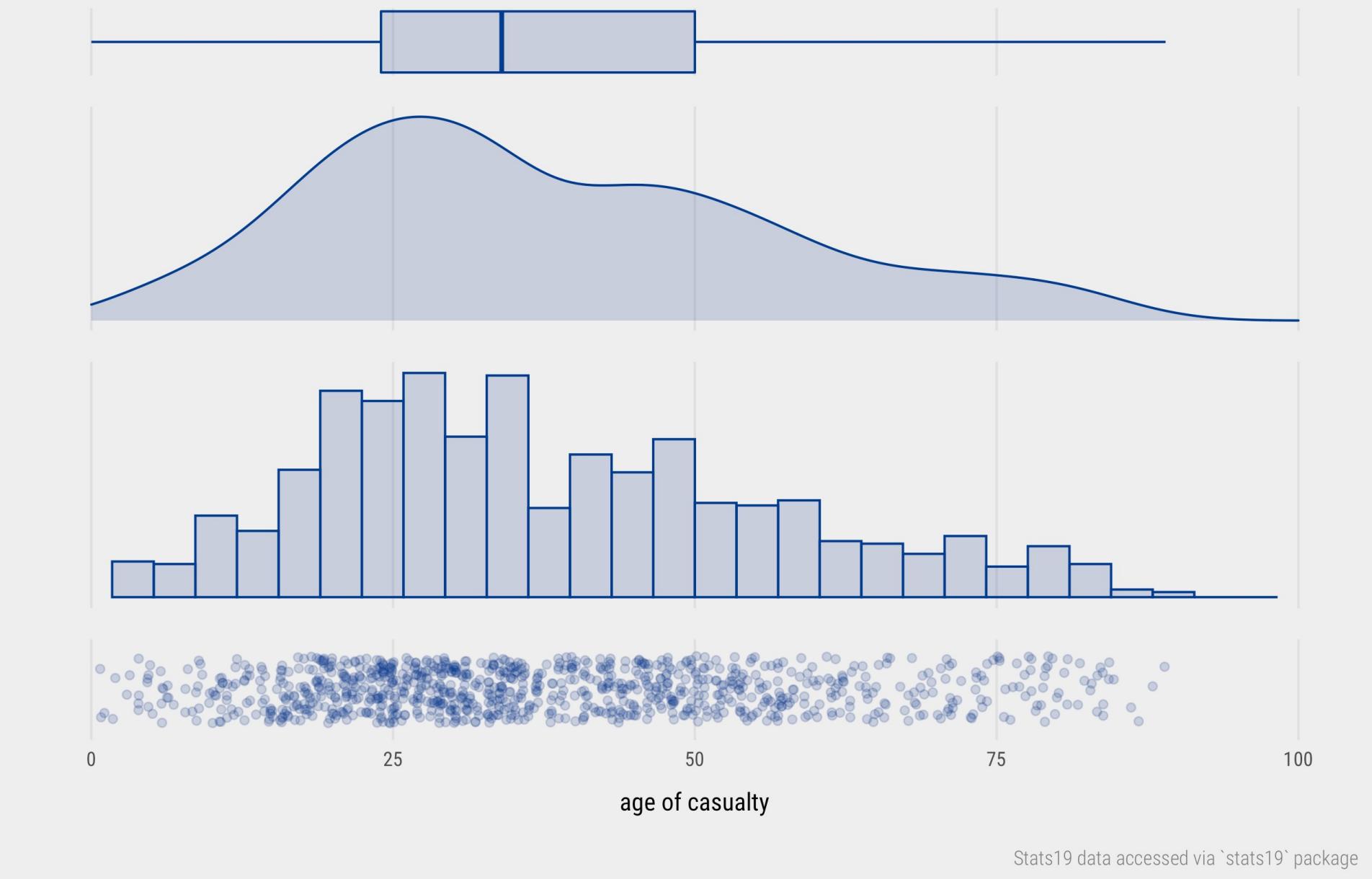
<https://www.r-graph-gallery.com/>



<https://seaborn.pydata.org/examples/index.html>

Plots of univariate distribution: age of casualty in Stats19 dataset

-- Strip-plot, histogram, density plot, boxplot | mean 36 years - median 33 years - mode 21 years



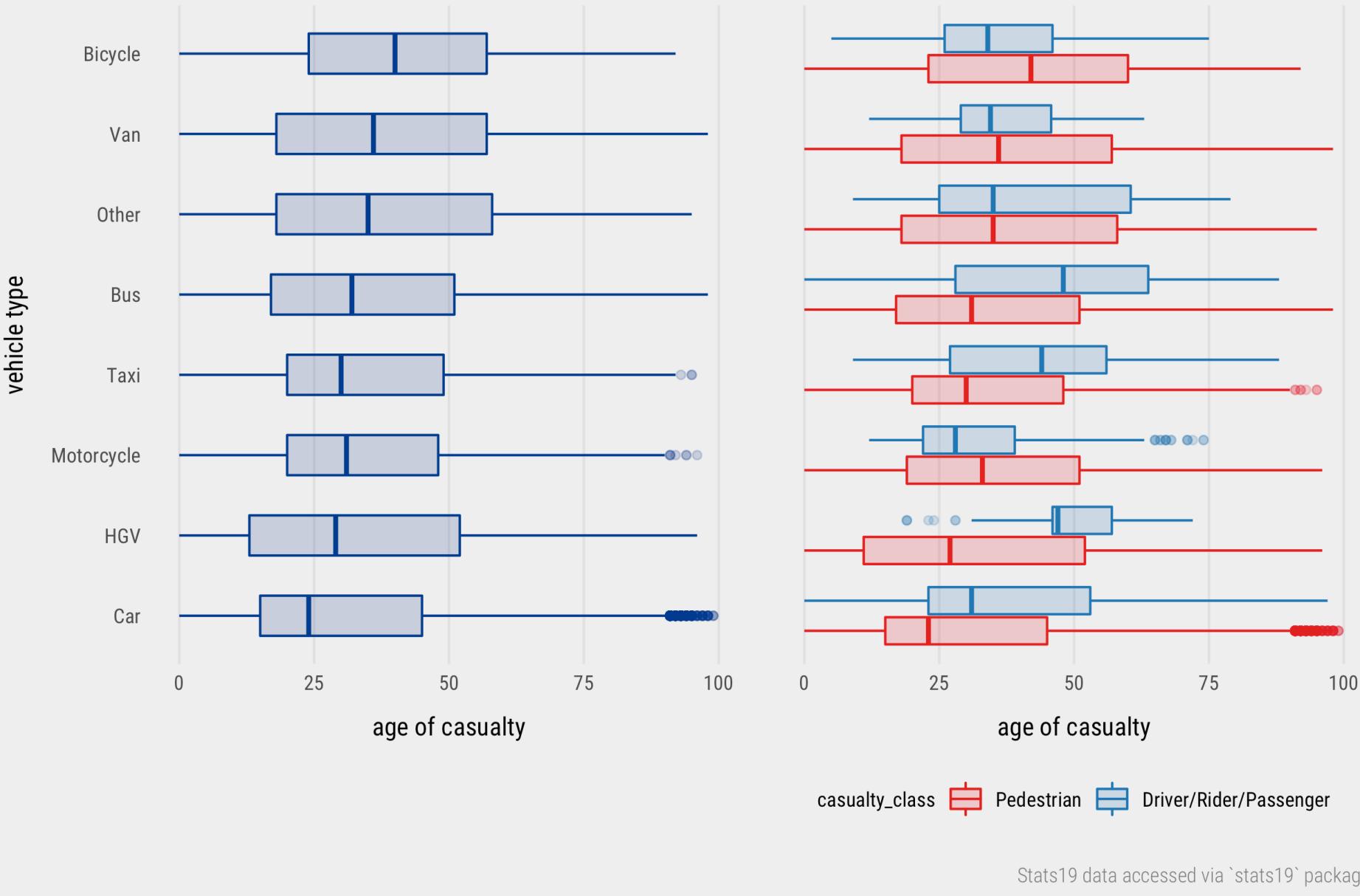
Stats19 data accessed via `stats19` package

Visualisation by Roger Beecham: <https://www.roger-beecham.com/vis-for-gds/class/04-class/>

Data: <https://cran.r-project.org/web/packages/stats19/vignettes/stats19.html>

Box plots of age of casualty by vehicle type, coloured by casualty class

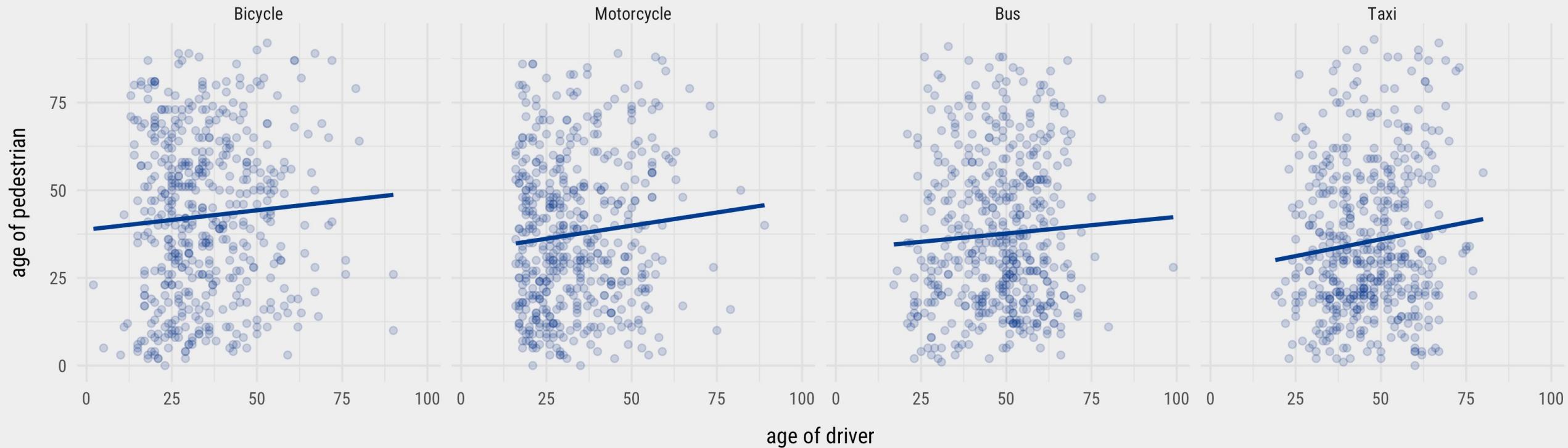
-- Random sample of 100k Stats19 Pedestrian-Vehicle crashes



Stats19 data accessed via `stats19` package

Scatterplots of pedestrian age by driver age and grouped by vehicle type

-- Random sample of Stats19 pedestrian-vehicle crashes stratified by vehicle type



Stats19 data accessed via `stats19` package

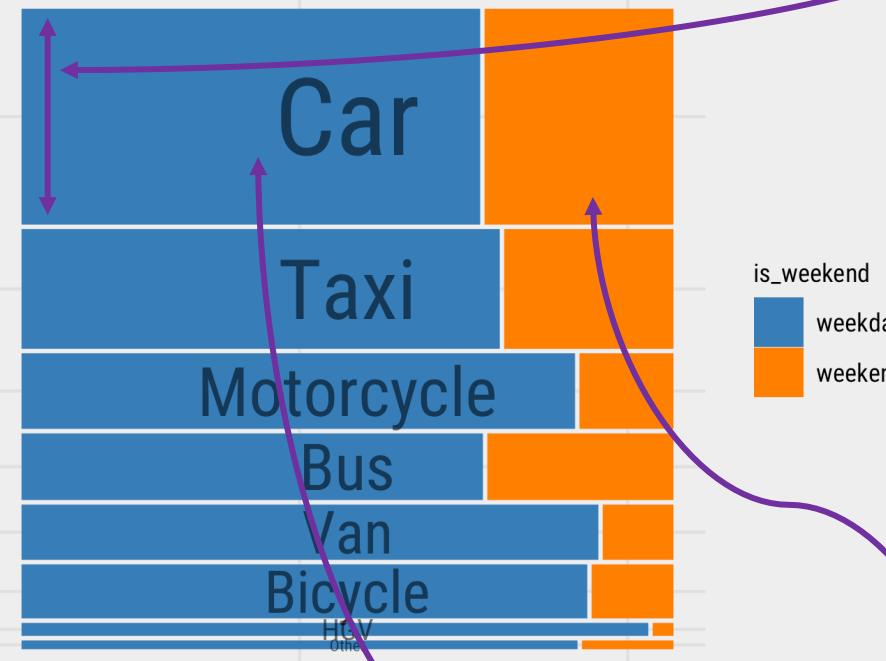
Mosaic Plot expld.

Proportion of incidents involving a **Car**

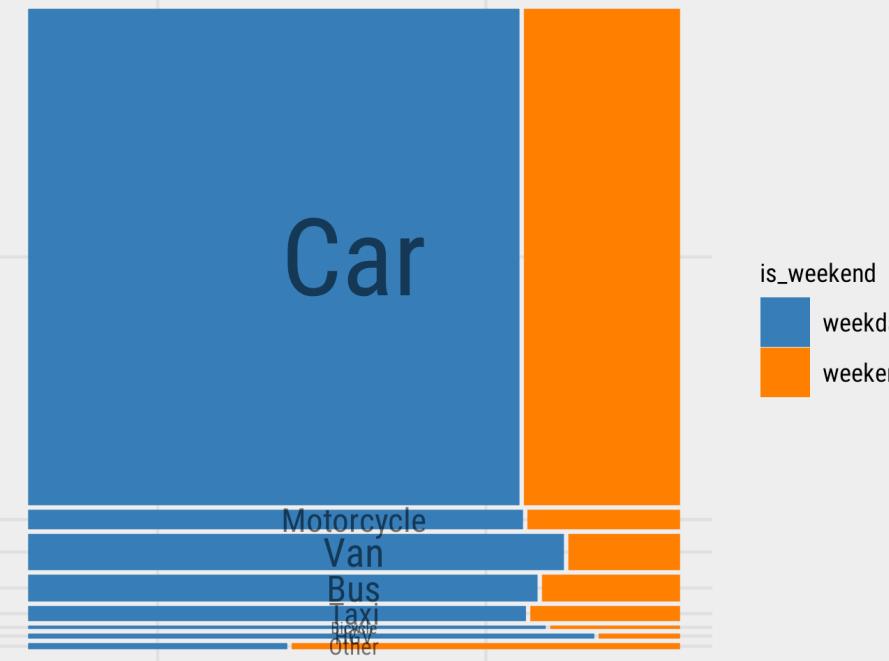
Pedestrian casualties by vehicle type and period in week

--Stats19 crashes 2010-2019

Westminster



Harrow



Proportion of **weekday** and **weekend**

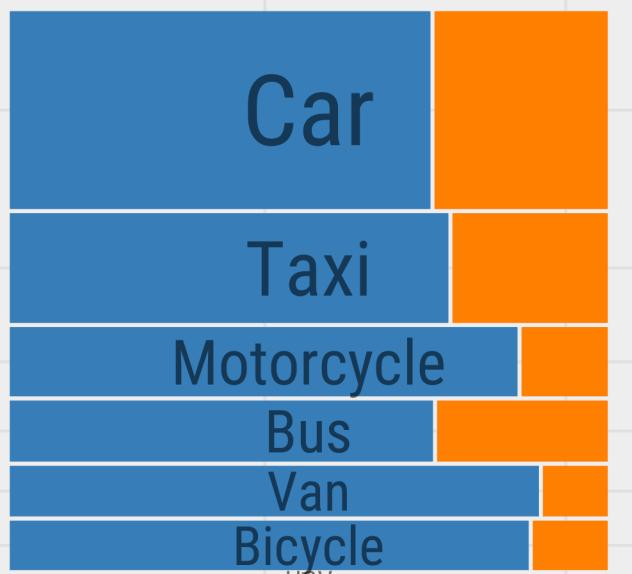
Stats19 data accessed via `stats19` package

RELATION – Vehicle Type vs. Time

Pedestrian casualties by vehicle type and period in week

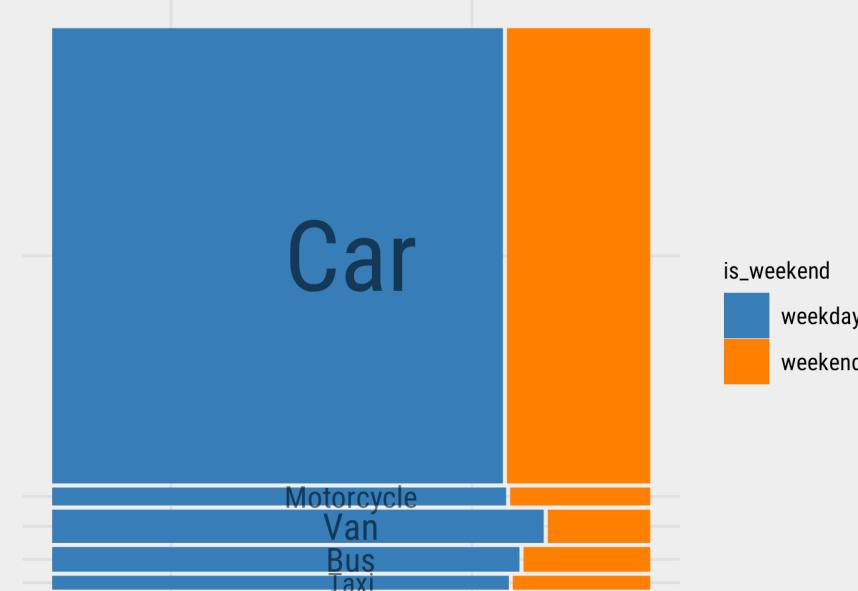
--Stats19 crashes 2010-2019

Westminster



COMPARISON (using two Mosaic plots)
– observing how the variations and relations vary under different conditions

Harrow



Stats19 data accessed via `stats19` package

Pedestrian casualties by vehicle type and period in week by London Borough

--Stats19 crashes 2010-2019



SMALL MULTIPLES

*Faceting & Sorting
(ordering)*

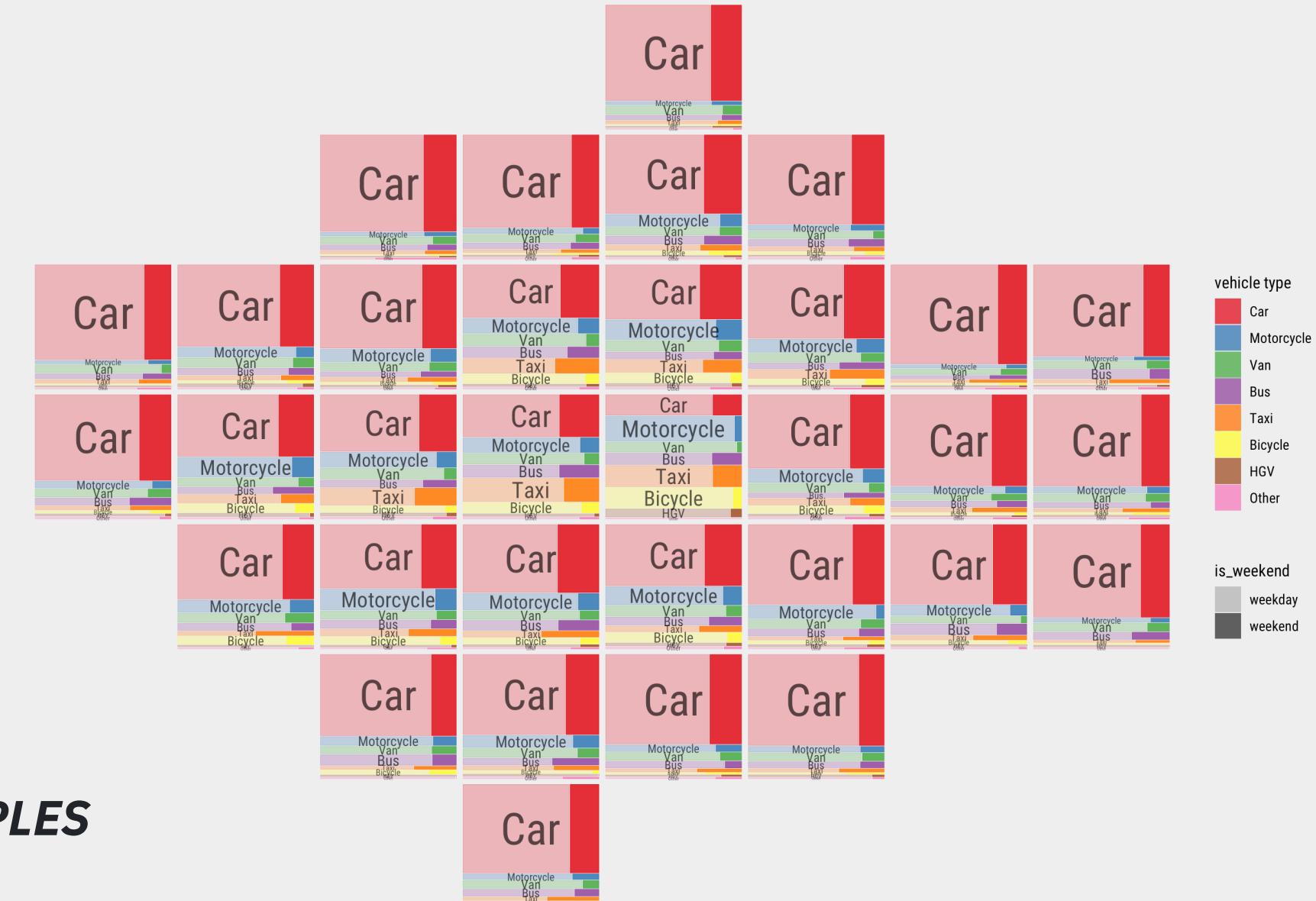
Data: <https://cran.r-project.org/web/packages/stats19/vignettes/stats19.html>

Visualisation by Roger Beecham:
<https://www.roger-beecham.com/vis-for-gds/class/04-class/>

Stats19 data accessed via 'stats19' package

Pedestrian casualties by vehicle type and period of week

--Relaxed spatial layout of London boroughs



SMALL MULTIPLES

Order semantics

Data: <https://cran.r-project.org/web/packages/stats19/vignettes/stats19.html>

Visualisation by Roger Beecham:
<https://www.roger-beecham.com/vis-for-gds/class/04-class/>

Stats19 data accessed via `stats19` package

COMPARISON –

Visual comparison for information visualization

Michael Gleicher¹, Danielle Albers¹, Rick Walker², Ilir Jusufi³,
Charles D. Hansen⁴ and Jonathan C. Roberts²

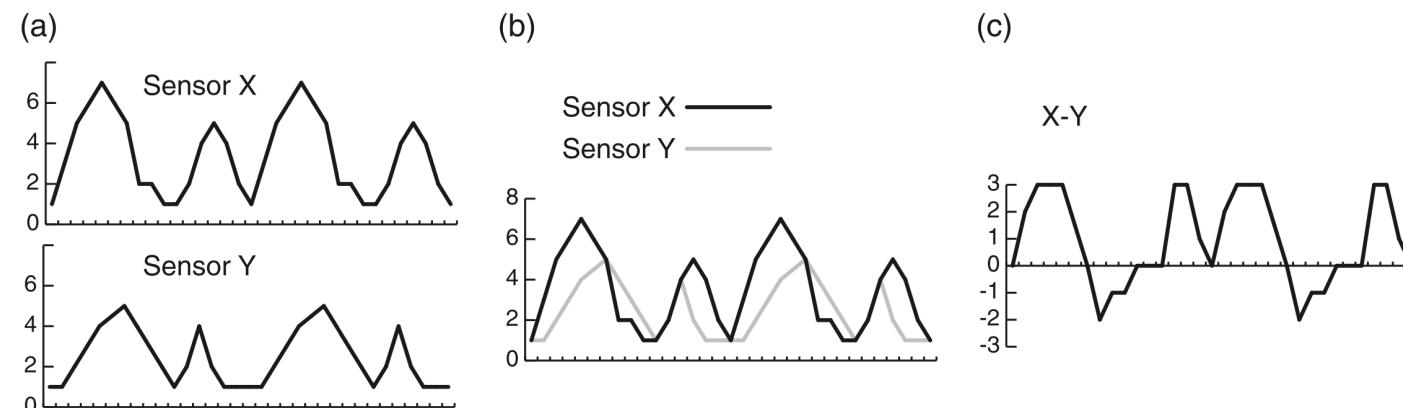
Information Visualization
10(4) 289–309
© The Author(s) 2011
Reprints and permissions:
sagepub.co.uk/journalsPermissions.nav
DOI: [10.1177/1473871611416549](https://doi.org/10.1177/1473871611416549)
ivi.sagepub.com
\$SAGE

Three broad design approach for comparison:

Juxtaposition designs place objects separately in either time or space.

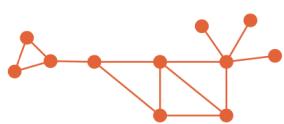
Superposition designs overlay multiple objects, presenting them at the same place and time.

Explicit encodings compute the relationships between objects and provide visual encoding of the relationships.





(a) Naïve Juxtaposition
(each network laid out independently)



(b) Juxtaposition (using
similar layouts to aid
comparison)



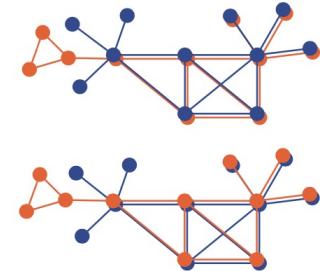
(c) Superposition (in the
same space)



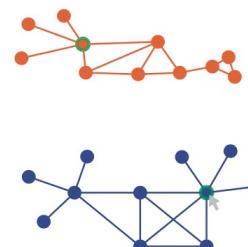
(d) Explicit Encoding:
Replacement (upper:
union graph, lower:
intersection graph)



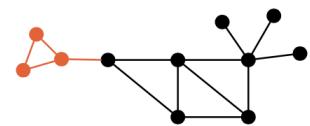
(e) Explicit Encoding:
Additive (members of
the intersection shown
added to one of the
graphs)



(f) Juxtaposition +
Superposition



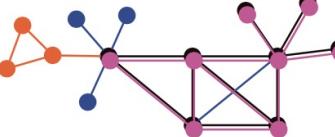
(g) Juxtaposition +
Explicit Encodings:
Using view coordination
to highlight
correspondences



(h) Juxtaposition +
Explicit Encoding:
Additive (intersections
added)



(i) Juxtaposition +
Explicit Encoding:
Abstraction (cliques
shown over juxtaposed
views)



(j) Superposition +
Explicit Encoding:
Overlay encoding



(k) Superposition +
Explicit Encoding: Ab-
straction+superposition
(cliques shown over
superimposed view)

A design space

For comparing two networks

COMPARISON - Layout and alignment!

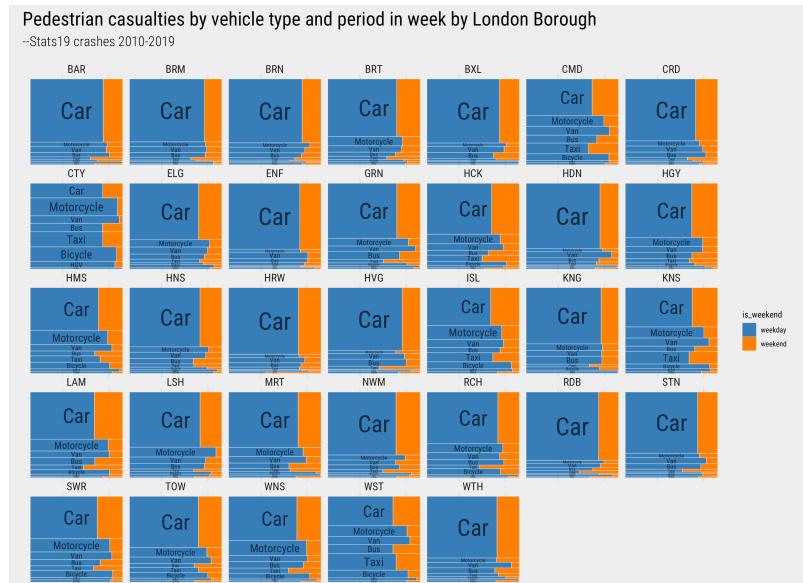
Aim for easy to read/compare individual graphics

Order is very important (everything?)

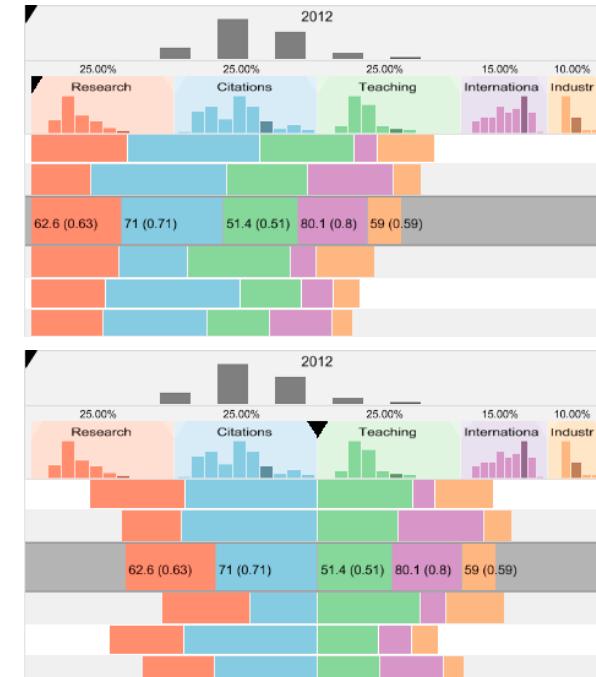
Time? Space? Attributes?

Label effectively

Align carefully



Visualisation by Roger Beecham: <https://www.roger-beecham.com/vis-for-gds/class/04-class/>

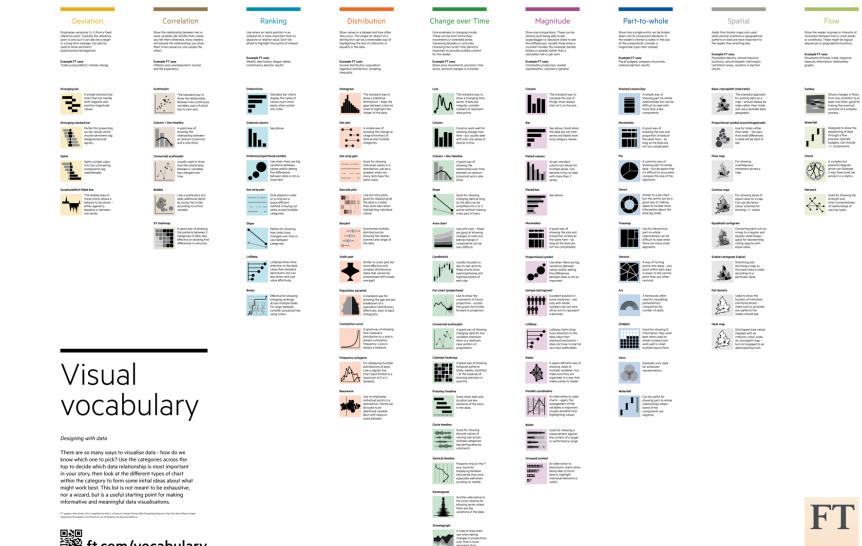


[LineUp: Visual Analysis of Multi-Attribute Rankings. Gratzl, Lex, Gehlenborg, Pfister, and Streit. IEEE Trans. Visualization and Computer Graphics (Proc. InfoVis 2013) 19:12 (2013), 2277–2286.]

P1: interrogate, relate, compare **variation, co-variation, deviation**

A few ideas to remember

- *Rich body of techniques, a lot of guidelines (e.g. FT vocabulary)*
- *Comparison -- a key strength of visualisation*
- *Facet and Sort – meaningfully and carefully*



Four key practices in Visual Data Science

... facilitated by interaction and visualisation

P1: interrogate, relate, compare **variation, co-variation, deviation**

P2: interact with data and computational tools

P3: generate **new perspectives** to “see” data and models

P4: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

What is Interaction for Data Visualization?

Evanthia Dimara and Charles Perin*

*“Interaction for visualization is the **interplay** between a person and a data interface involving a **data-related intent**, at least one **action** from the person and an interface **reaction** that is perceived as such”*

Toward a Deeper Understanding of the Role of Interaction in Information Visualization

Ji Soo Yi, Youn ah Kang, John T. Stasko, *Member, IEEE*, and Julie A. Jacko

“.. the concept of ‘What a user wants to achieve’, herein described as “**user intent**”, is quite effectively classified as these **low-level interaction techniques** : ”

Select: mark something as interesting

Explore: show me something else

Reconfigure: show me a different arrangement

Encode: show me a different representation

Abstract/Elaborate: show me more or less detail

Filter: show me something conditionally

Connect: show me related items

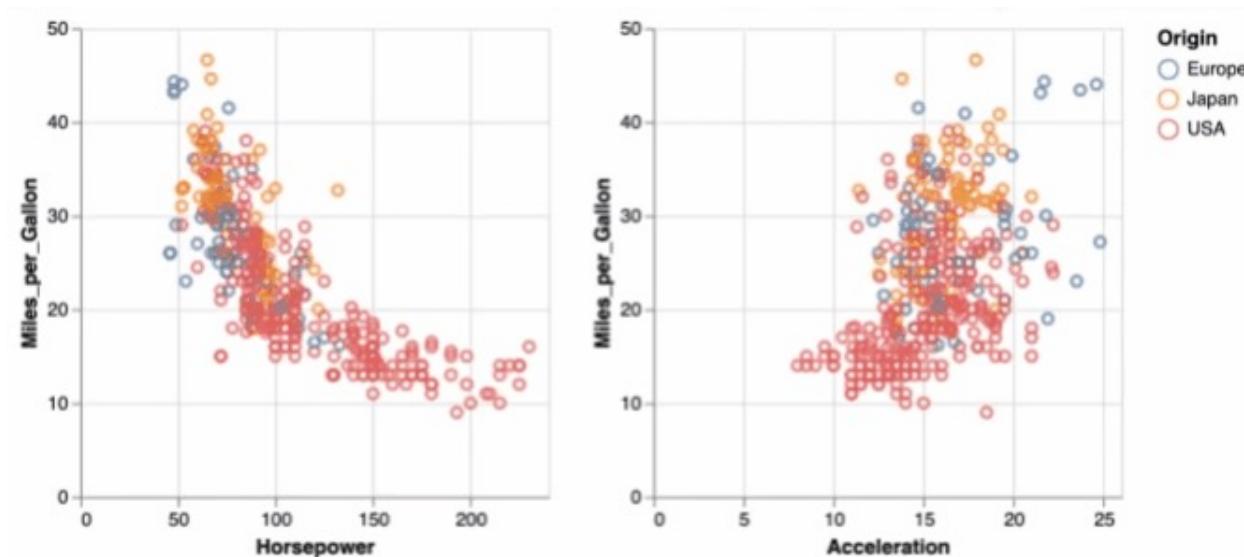
Linked highlighting, aka. brushing and linking

explore RELATIONs between variables

To see how data contiguous in one view are distributed within another

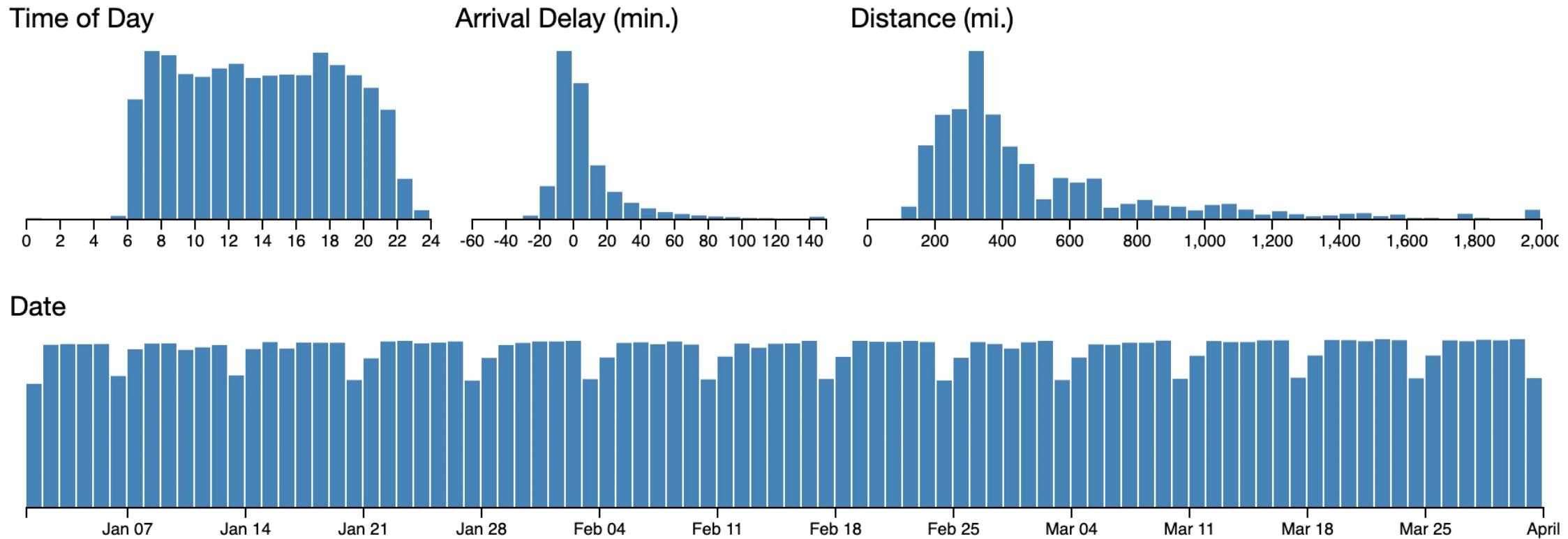
Fundamental interaction technique within interactive EDA

Focus + Context



cross filtering

- Focused on filtering items and narrowing down to subsets
- coordinated views/controls combined
 - update when any ranges change



March 31, 2001

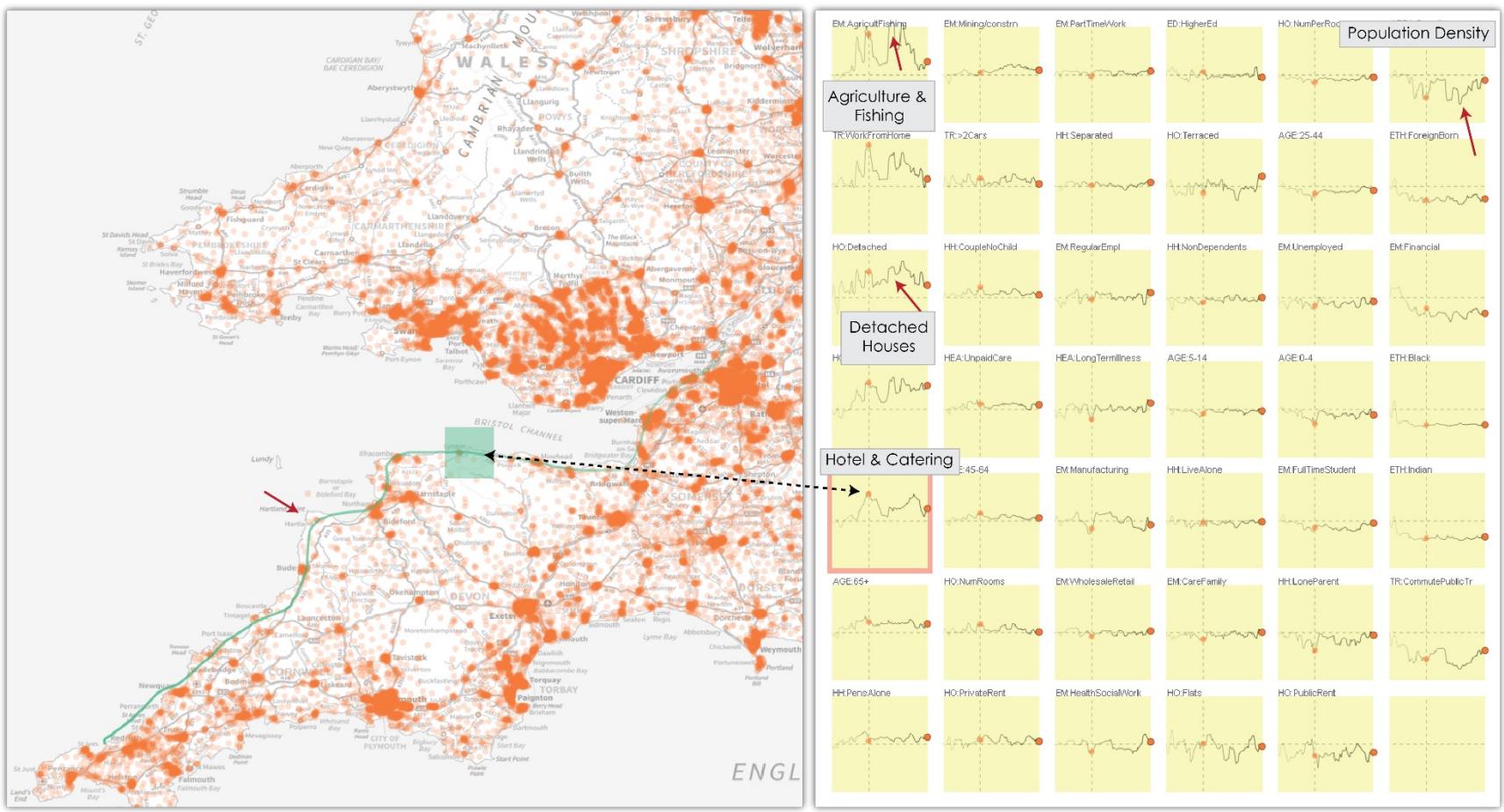
231,083 of 231,083 flights selected.

[\[http://square.github.io/crossfilter/\]](http://square.github.io/crossfilter/)

STUDY 1 – MULTIVARIATE GEOGRAPHICAL DATA

Attribute Signatures:
Dynamic Visual
Summaries for
Analyzing Multivariate
Geographical Data

*Cagatay Turkay, Aidan Slingsby,
Helwig Hauser, Jo Wood, Jason
Dykes, InfoVis 2014*



UK Census of Population in 2001 and 2011 for the
181,000 Output Areas (OA)
for 41 indicators



CITY UNIVERSITY
LONDON

question is...

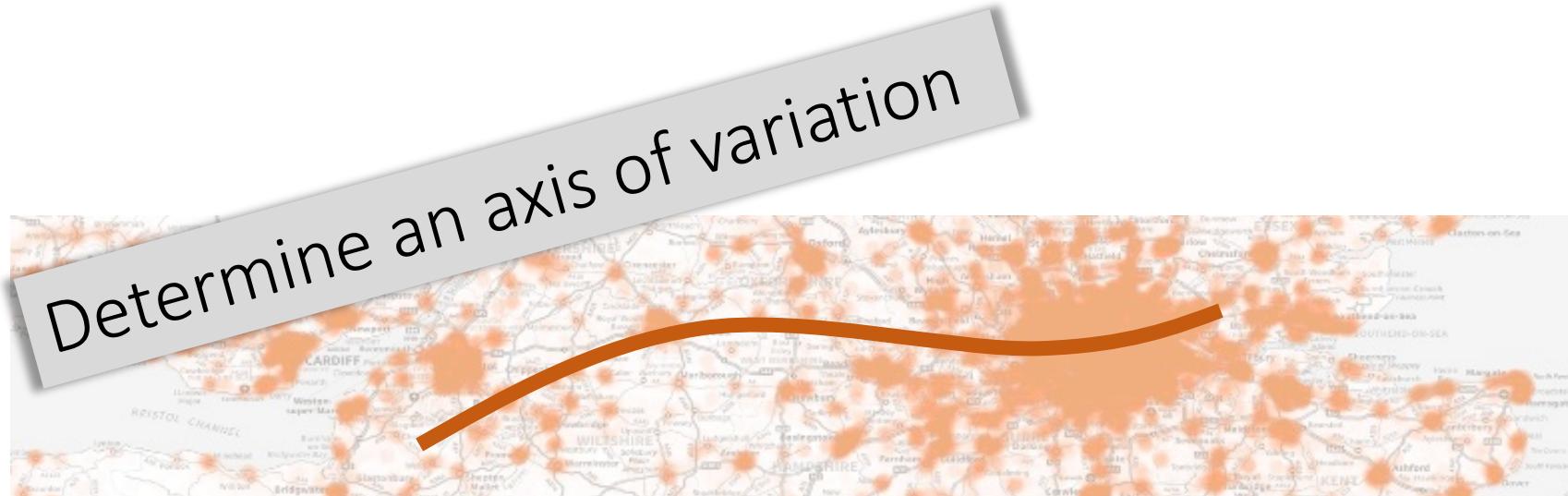
How all the variables
vary over space (and time) ?

we built methods to:

interactively
generate
visual summaries of change
in all the variables
in response to variation
(location, extent, resolution)

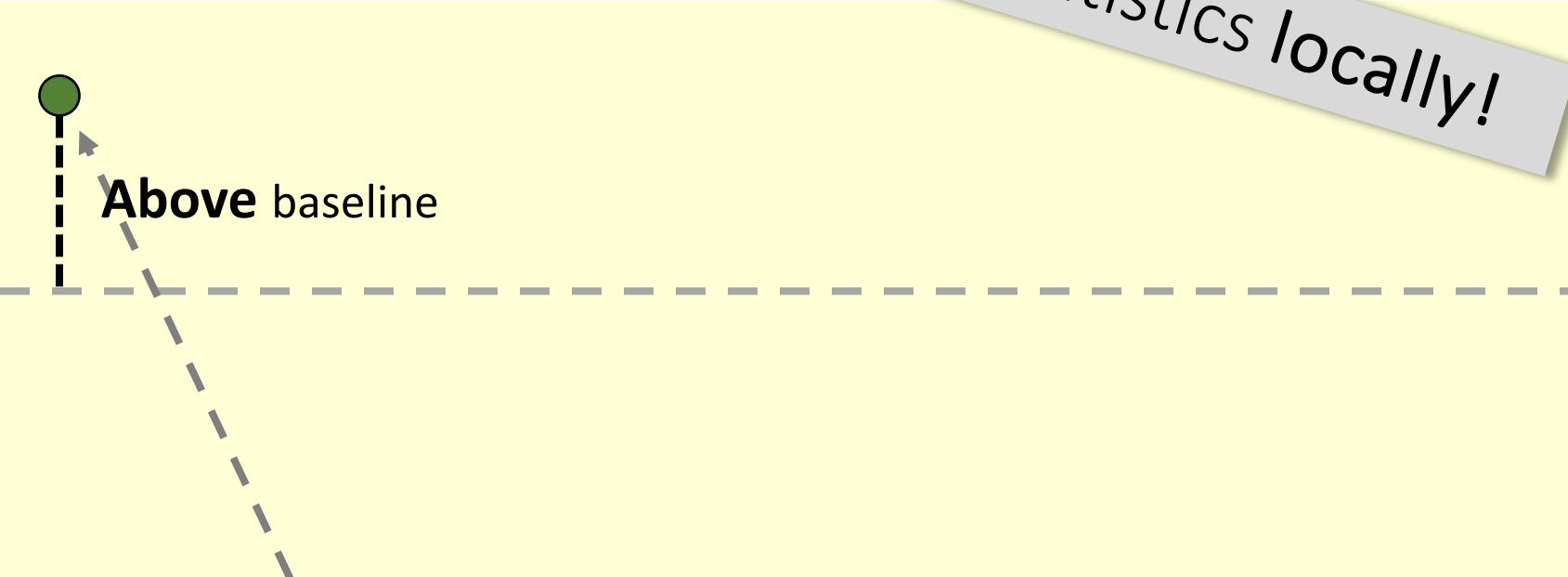
Attribute signatures

It starts with a **map** ...
and **an attribute** to analyze,
e.g., unemployment rates



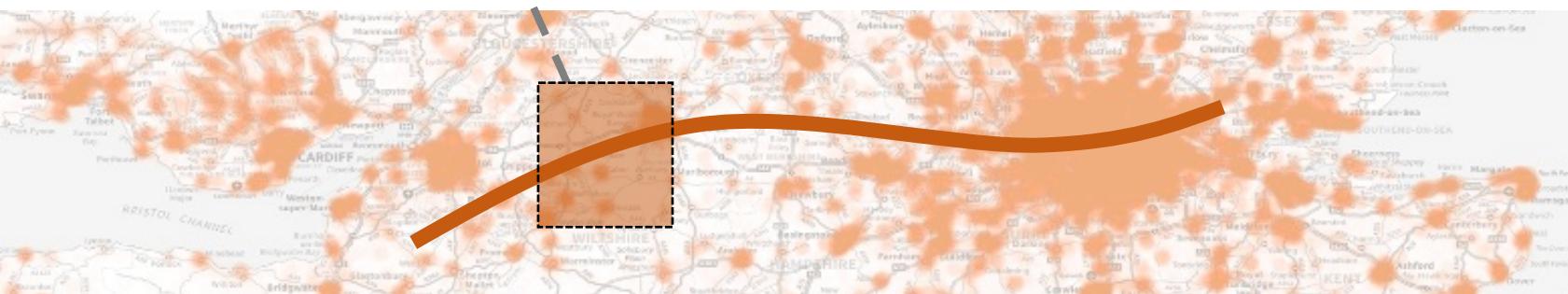
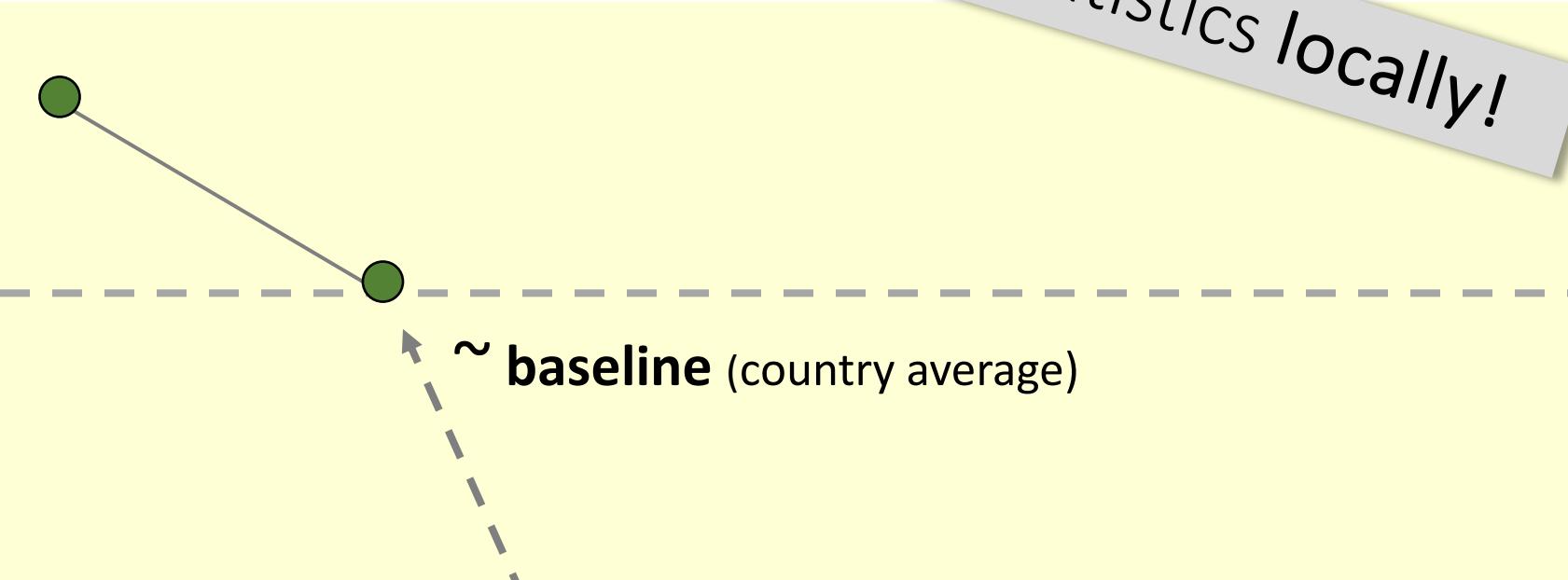
Attribute signatures

Compute statistics locally!



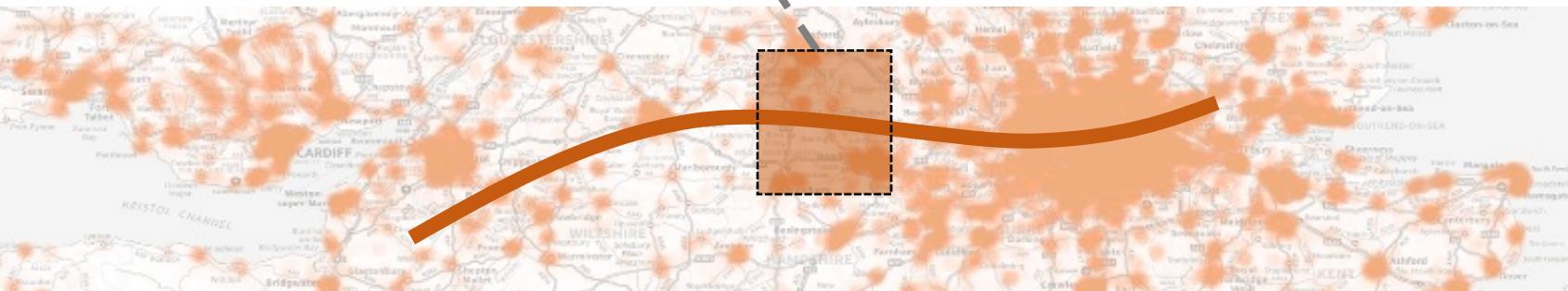
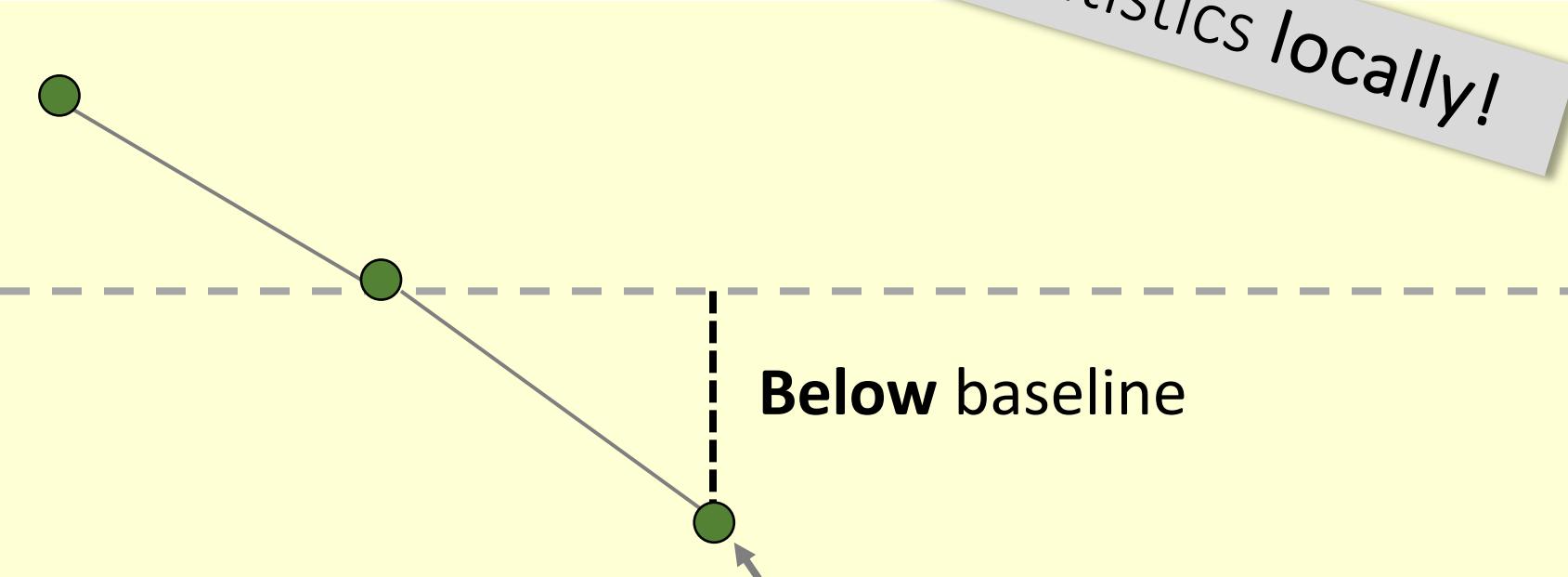
Attribute signatures

Compute statistics locally!



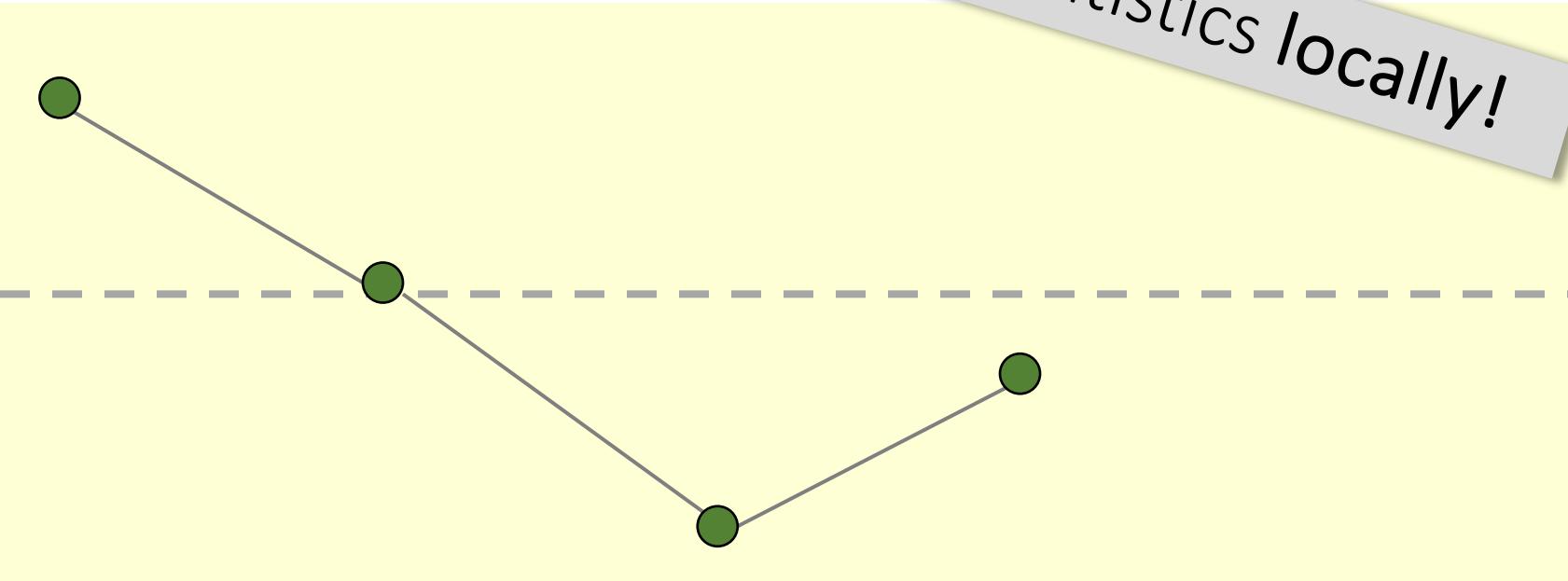
Attribute signatures

Compute statistics locally!



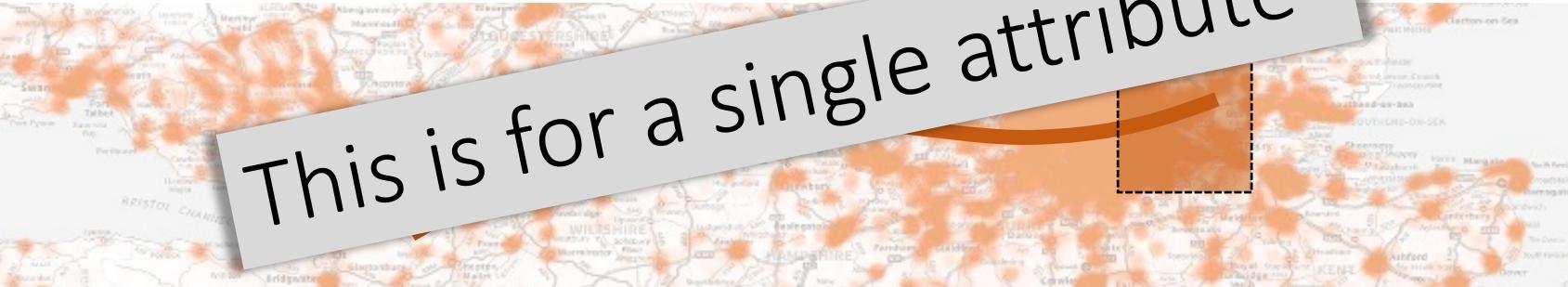
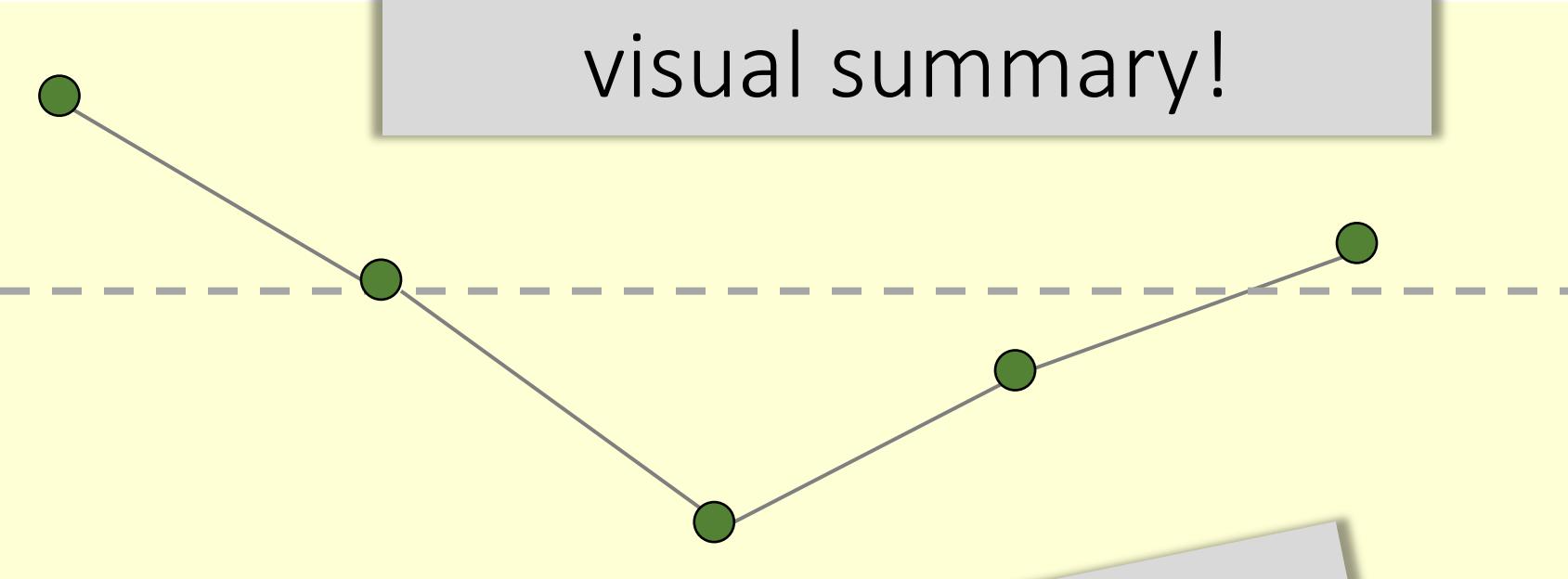
Attribute signatures

Compute statistics locally!



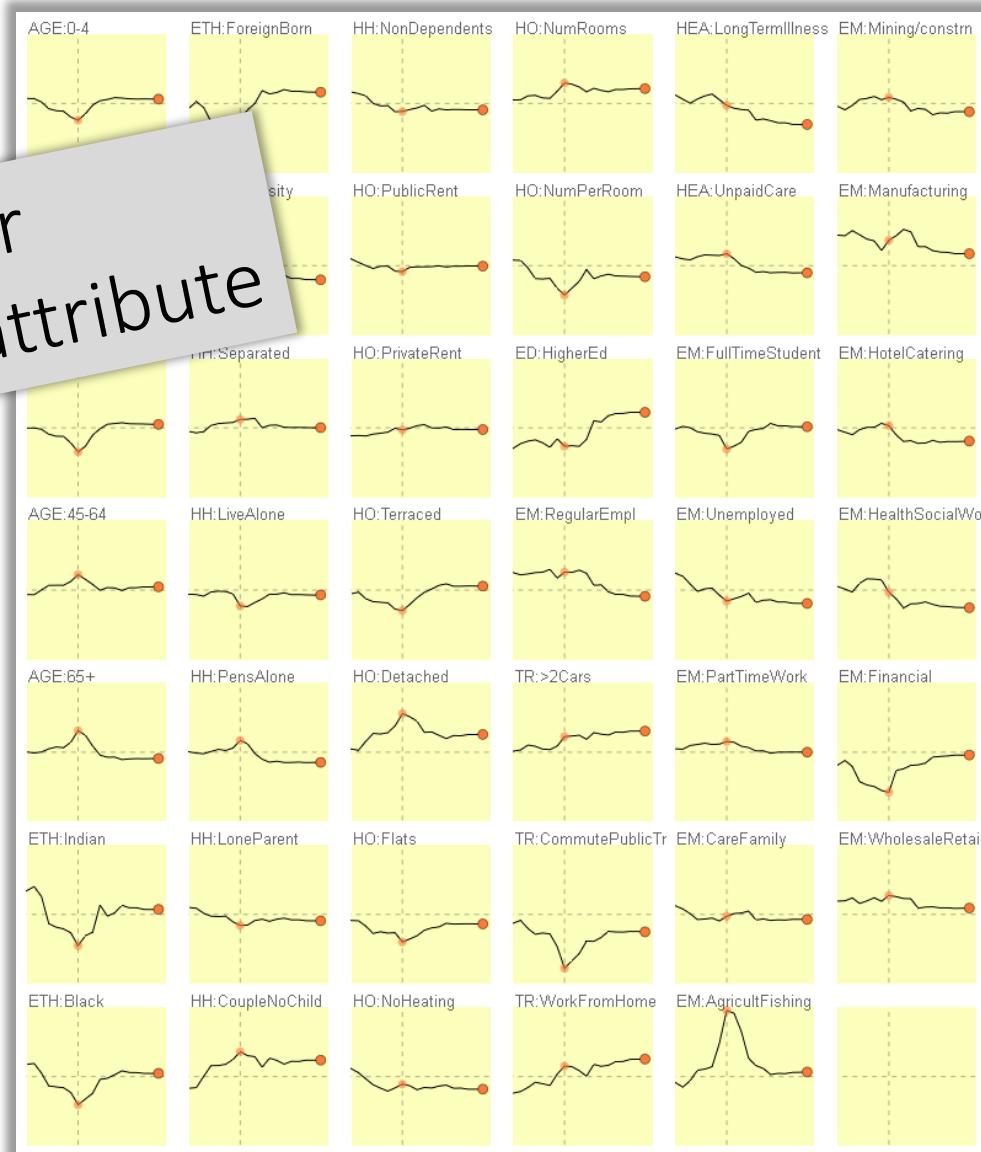
Attribute signatures

A dynamically generated
visual summary!

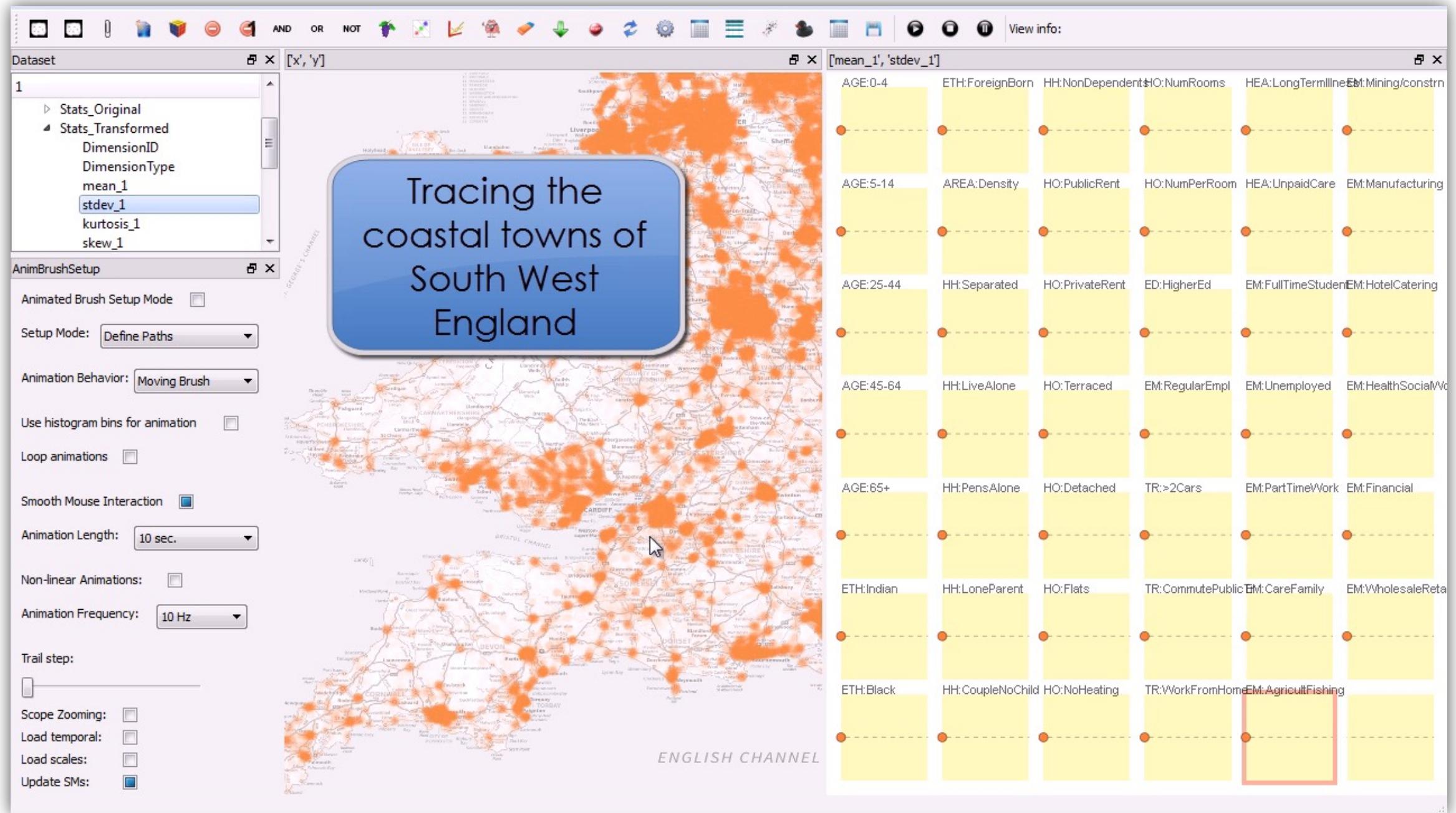


How about several attributes?

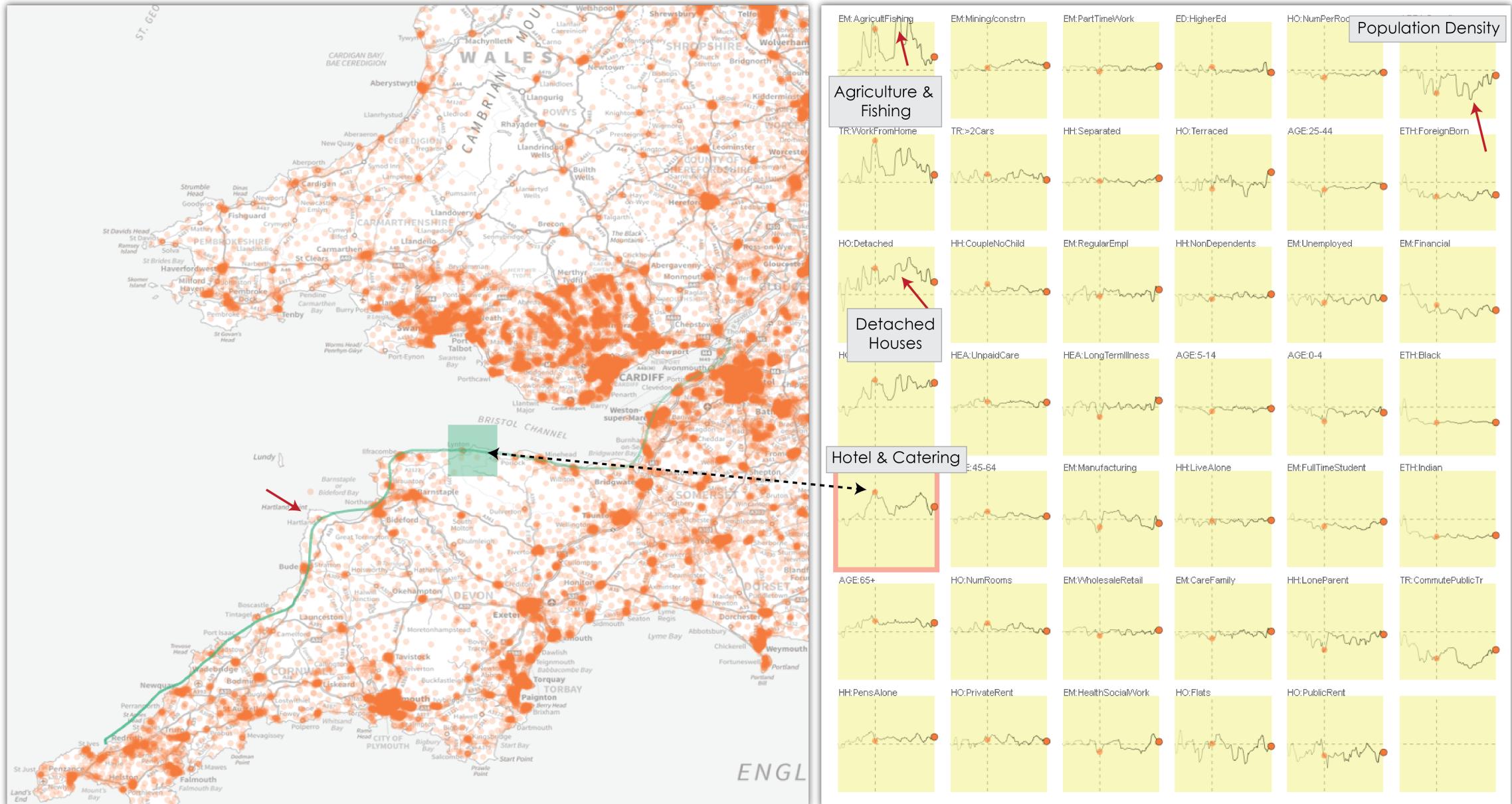
One for
each attribute



Linked
small
multiples



Along the Southwest England coast



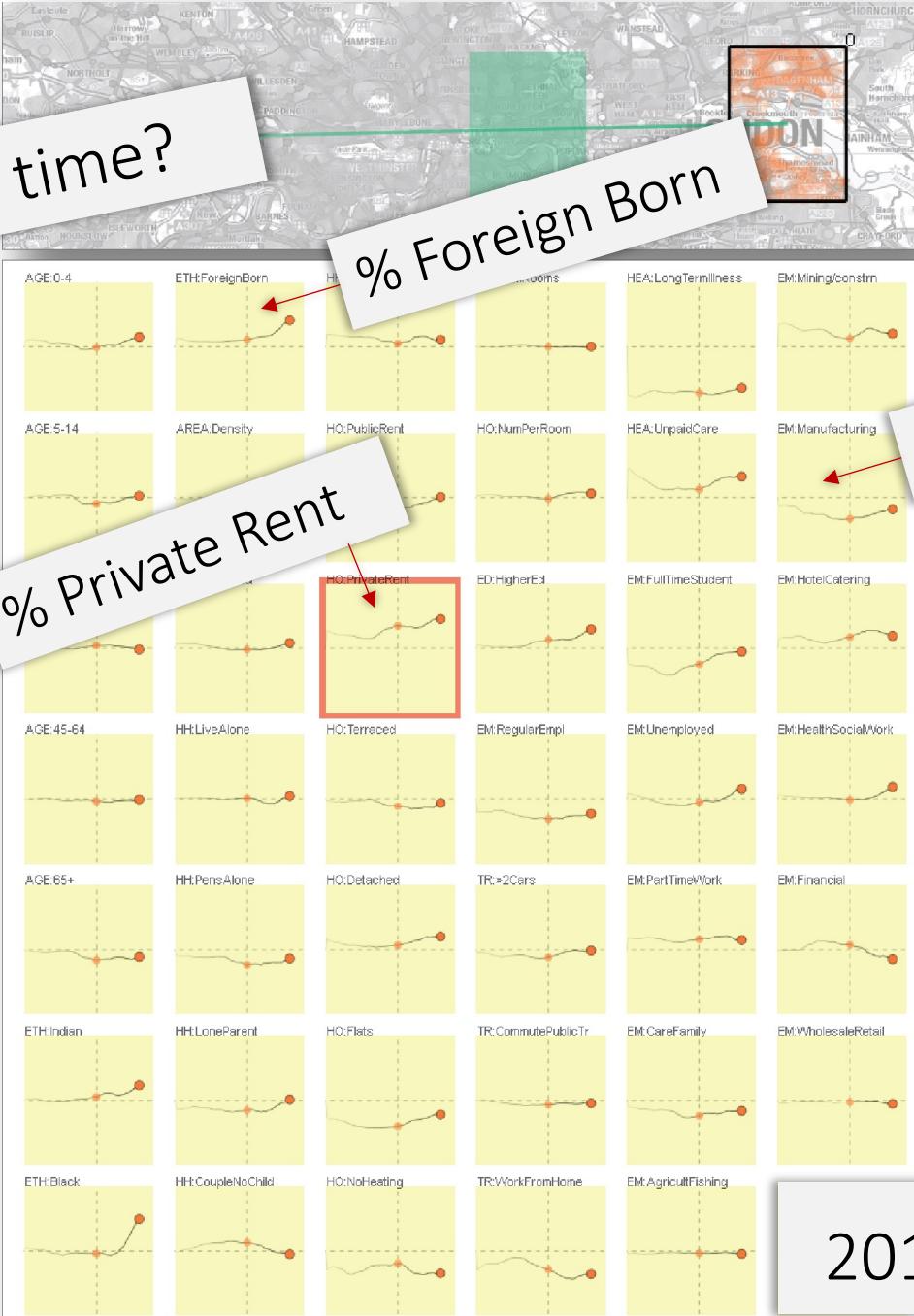
How about time?

% Foreign Born

% Manufacturing

% Private Rent

2011 vs. 2001



P2: interact with data and computational tools

A few ideas to remember

- *Interaction as an enabler of the VDS dialogue*
 - *for the analyst to “express” tacit/domain knowledge, e.g., selecting an area of interest, leaving out values..*
 - *for the machine to “understand” (and learn)*
- *Deal with complexity in working with data, e.g., filtering to a subset, narrowing the focus to a few variables*
- *A cornerstone in VDS to iteratively generate and evaluate hypotheses at speed*

Four key practices in Visual Data Science

... facilitated by interaction and visualisation

P1: interrogate, relate, compare **variation, co-variation, deviation**

P2: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

P3: generate **new perspectives** to “see” data and models

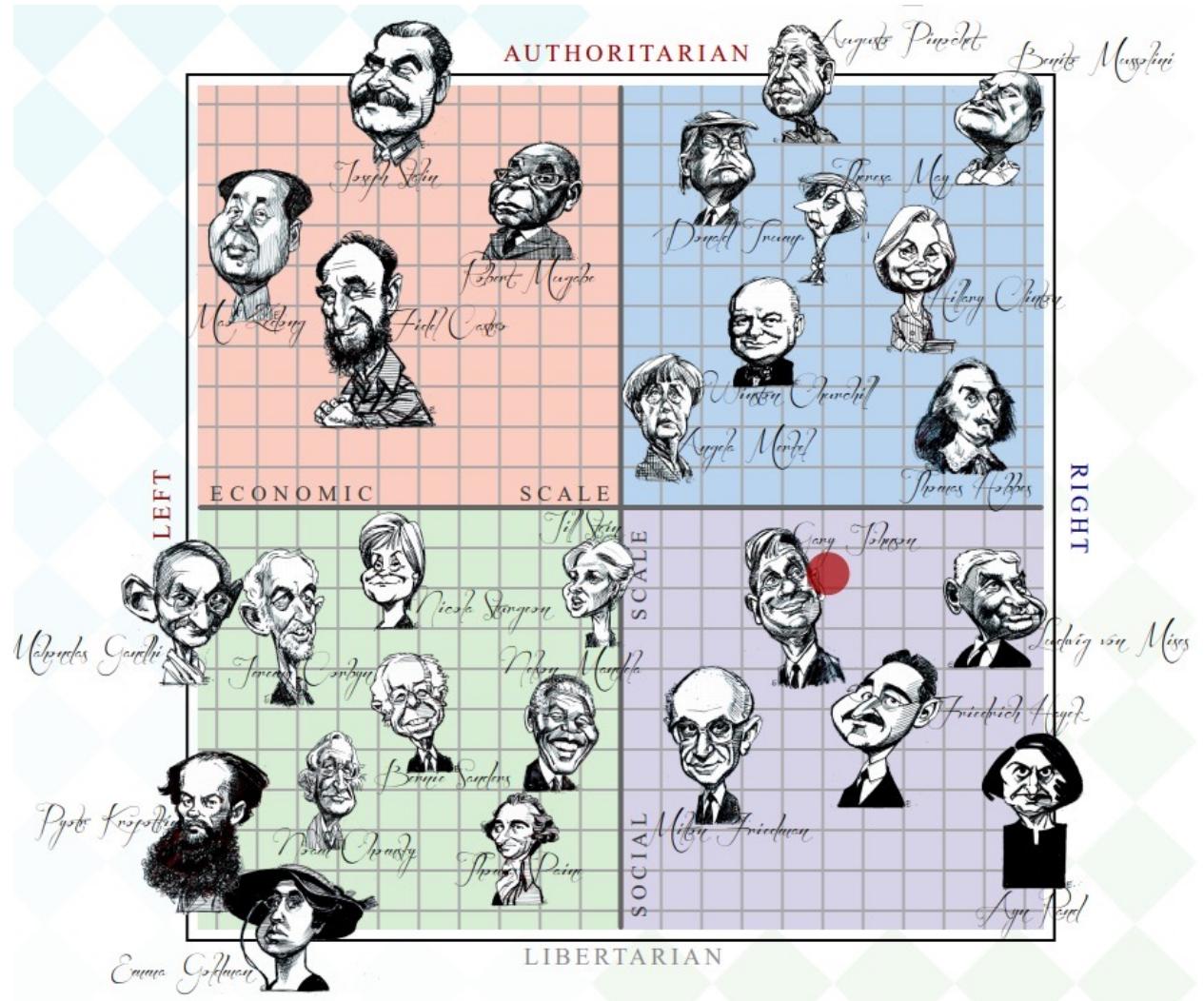
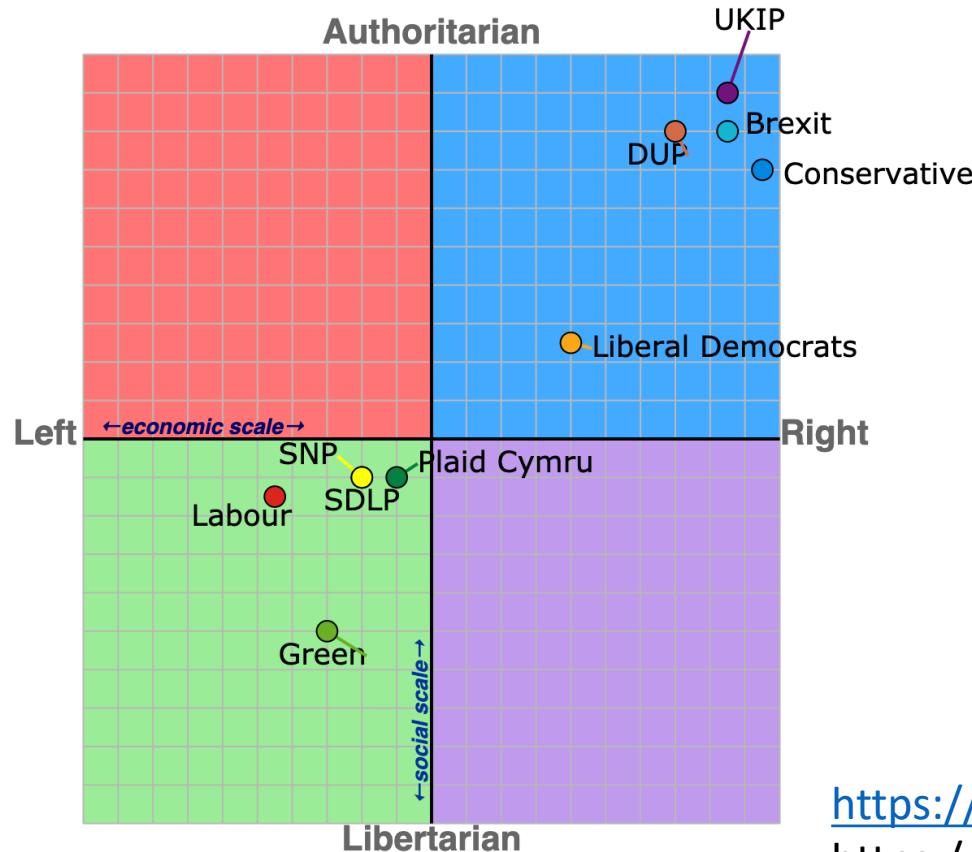
P4: interact with data and computational tools





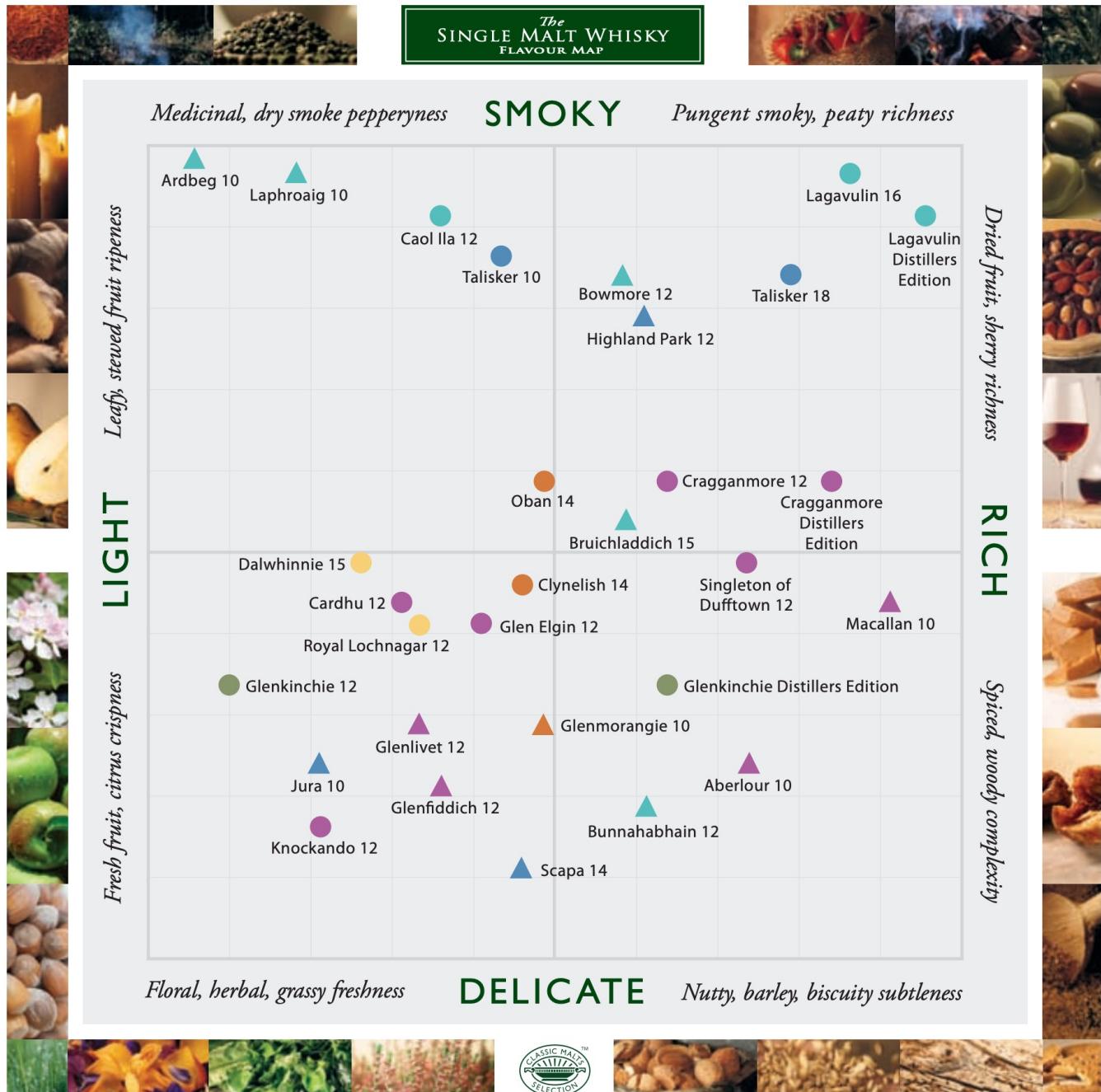
UK Parties 2019 General Election

7 December 2019



<https://www.politicalcompass.org/uk2019>

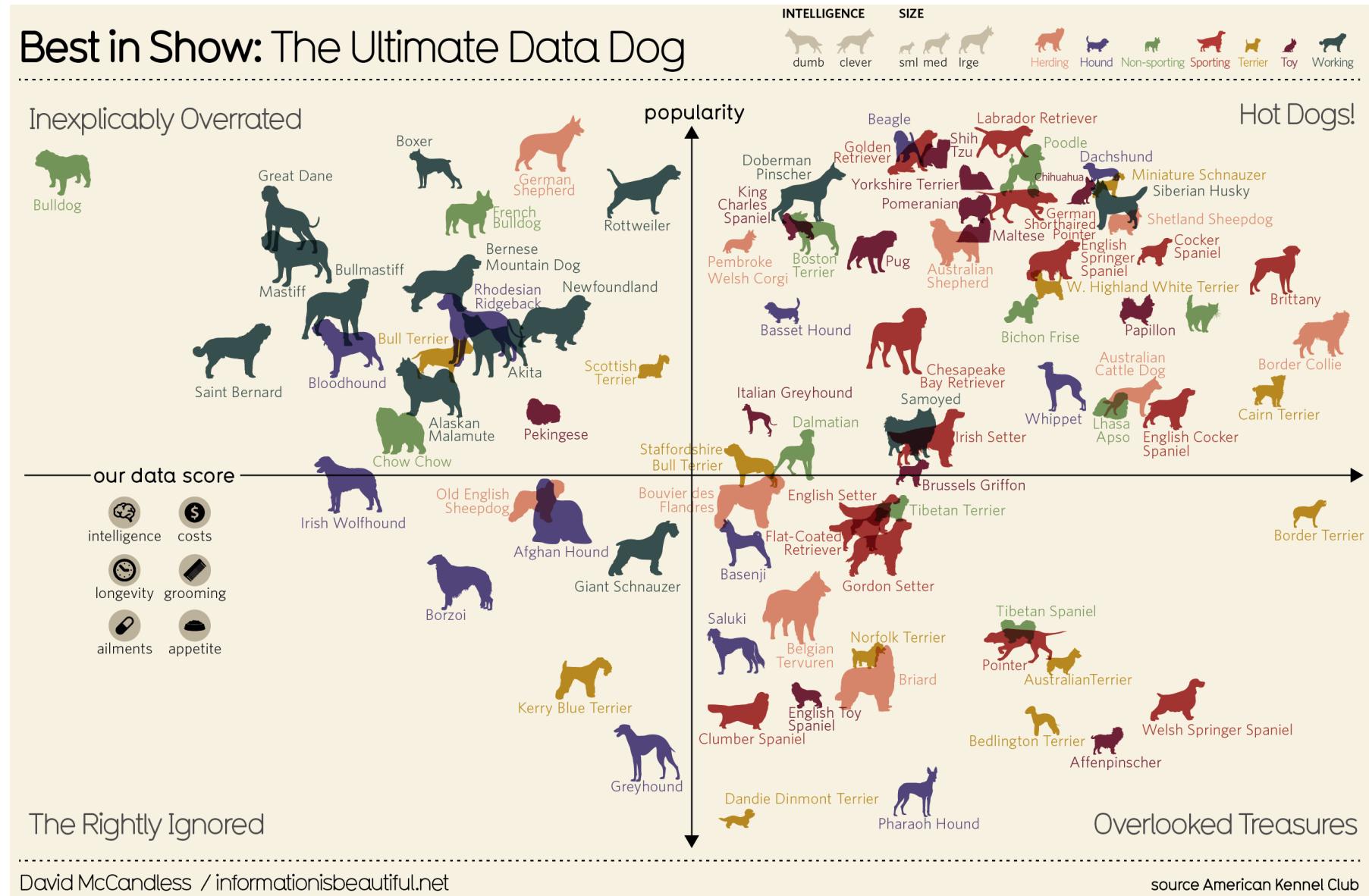
<https://danieljmitchell.wordpress.com/2018/05/16/the-political-compass-test/>



Malts with this shape symbol are part of the CLASSIC MALTS Selection owned by Diageo Scotland Limited.
The colours relate to the regional map of Scotland on the inside cover.

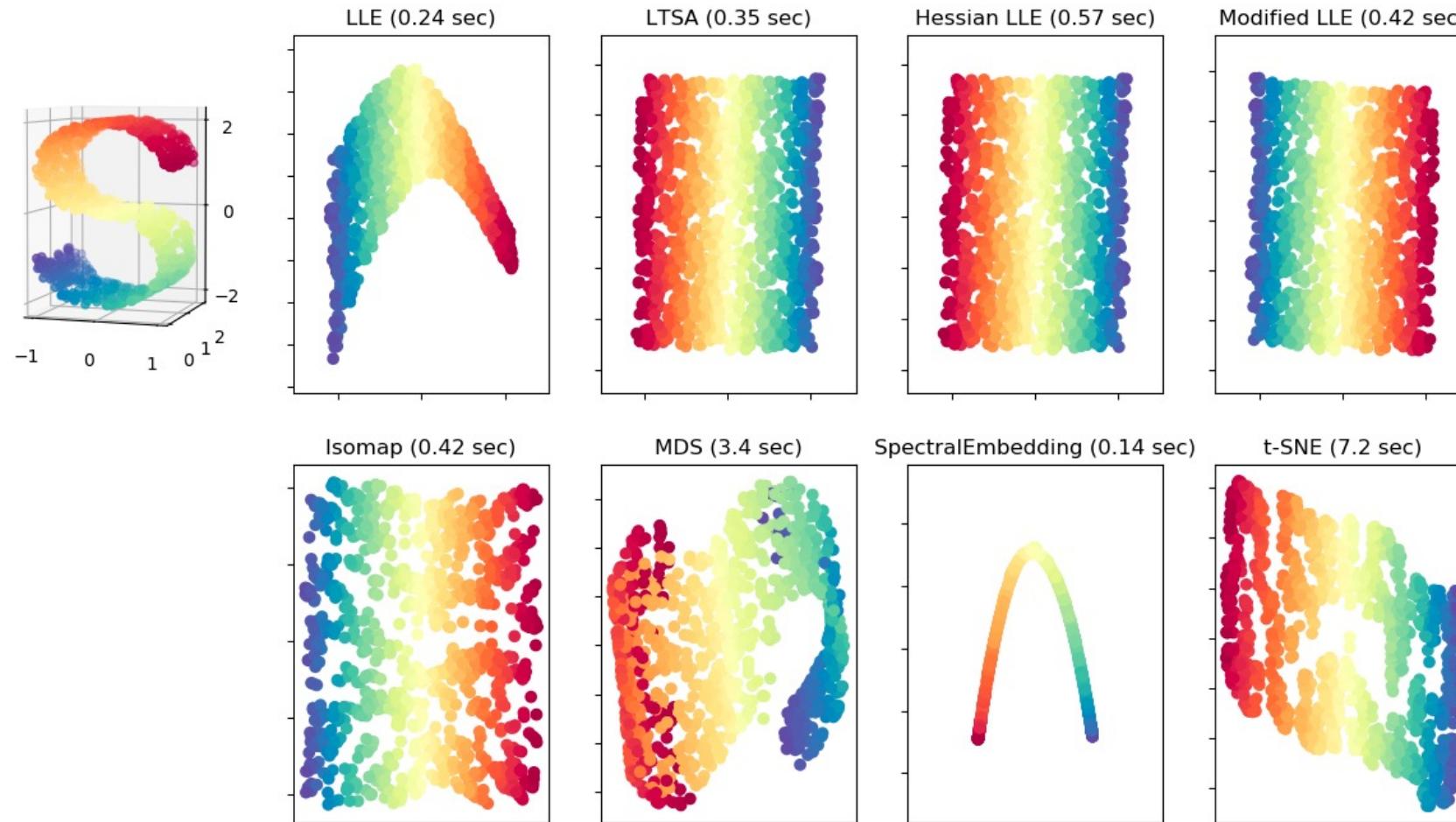
Malts with this shape symbol are not part of the CLASSIC MALTS Selection
<https://www.scotchkmaltwhisky.co.uk/whiskyflavourmap.pdf>

Feature generation to provide new “axes” for data analysis



Algorithms that provide new “spaces”, think of PCA, MDS, t-SNE, uMAP...

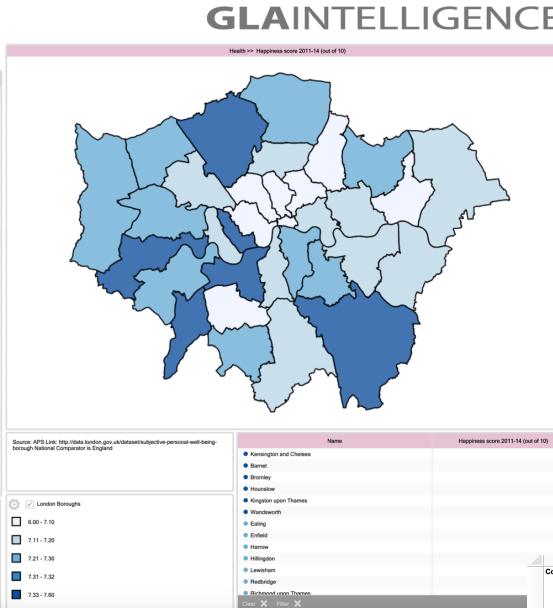
Manifold Learning with 1000 points, 10 neighbors



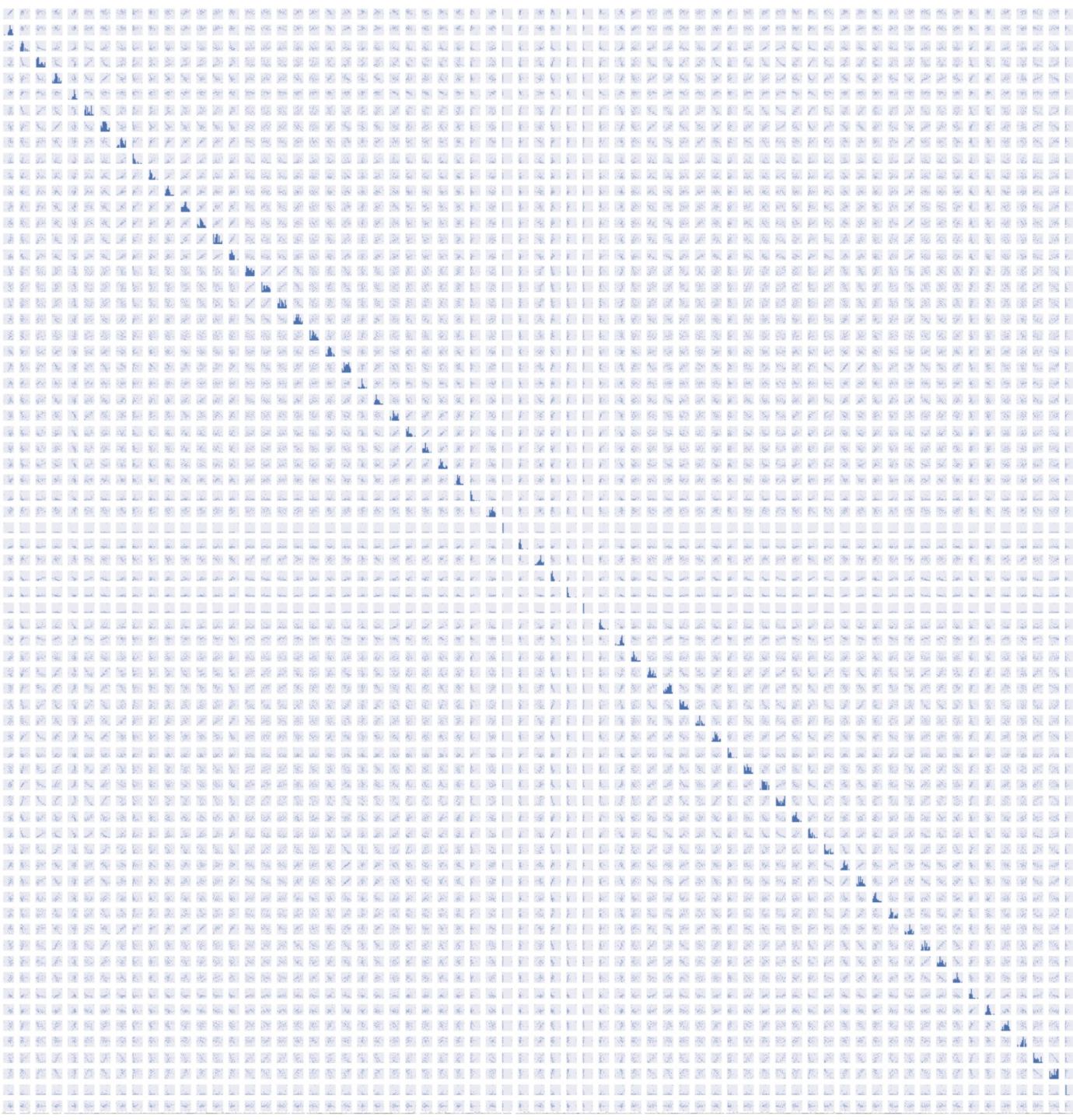
<http://scikit-learn.org/stable/modules/manifold.html>

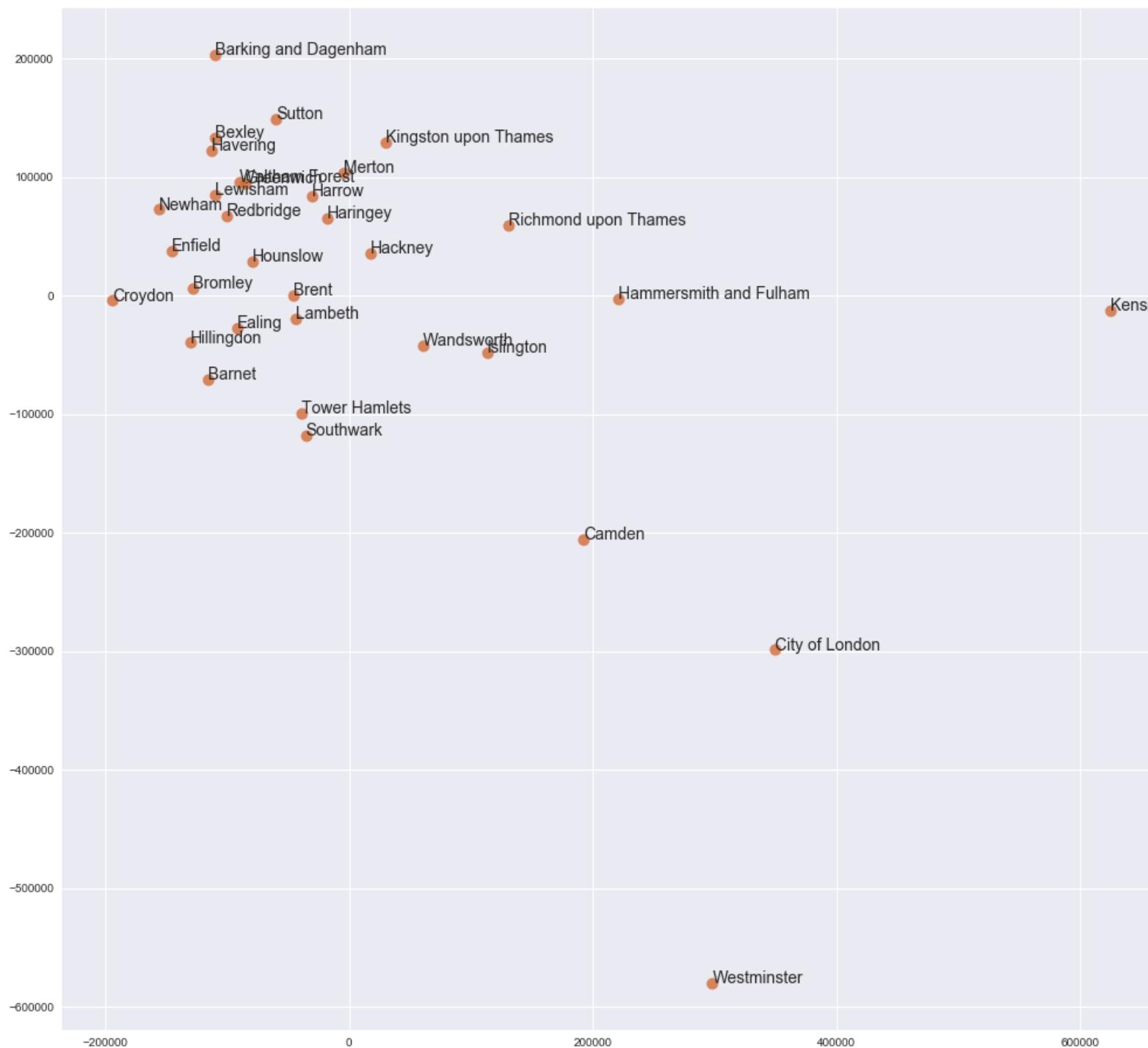
High-dimensional data – an example

LONDON BOROUGH PROFILES



" The London Borough Profiles help paint a general picture of an area by presenting a range of headline indicator data in both spreadsheet and map form to help show statistics covering demographic, economic, social and environmental datasets for each borough, alongside relevant comparator areas..."

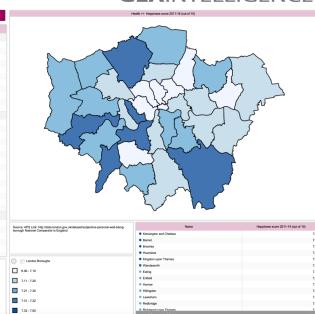




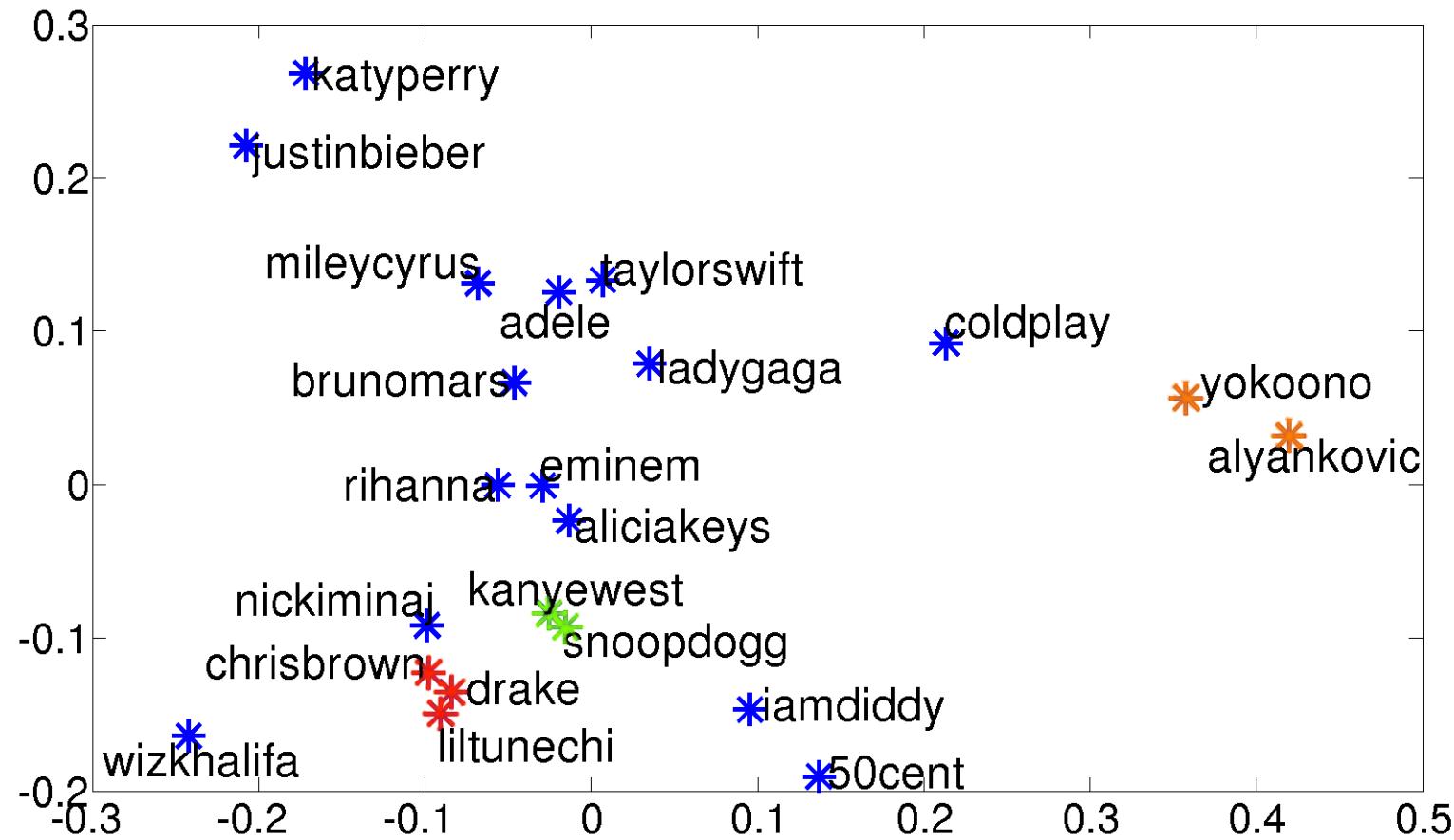
LONDON BOROUGH PROFILES

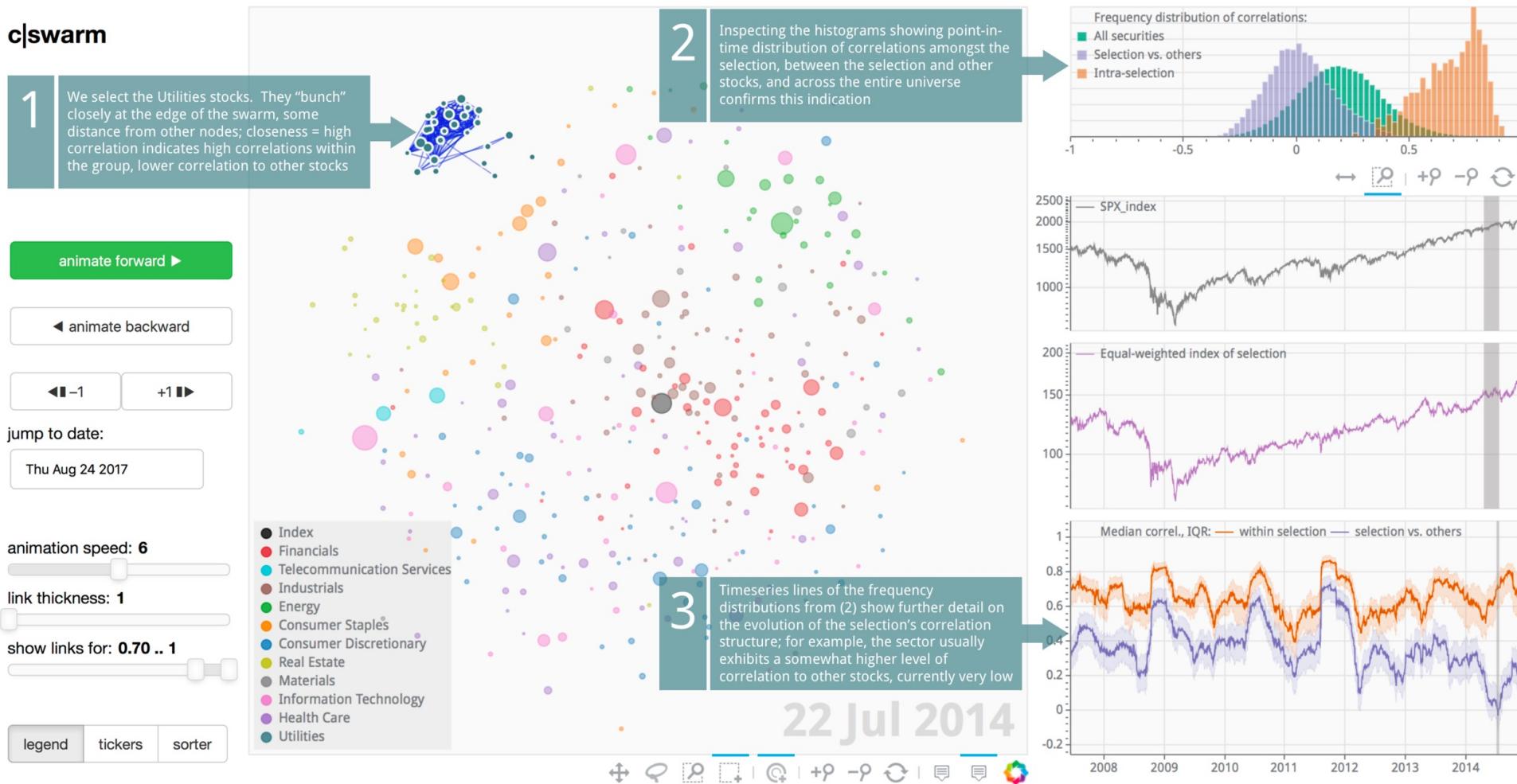
Demographic Data	Health	Environment & Land	Community
Demographic Data	Health	Environment & Land	Community

GLAINTELLIGENCE



Distance semantics matter (e.g., co-following on Twitter)





Hunting High and Low: Visualising Shifting Correlations in Financial Markets

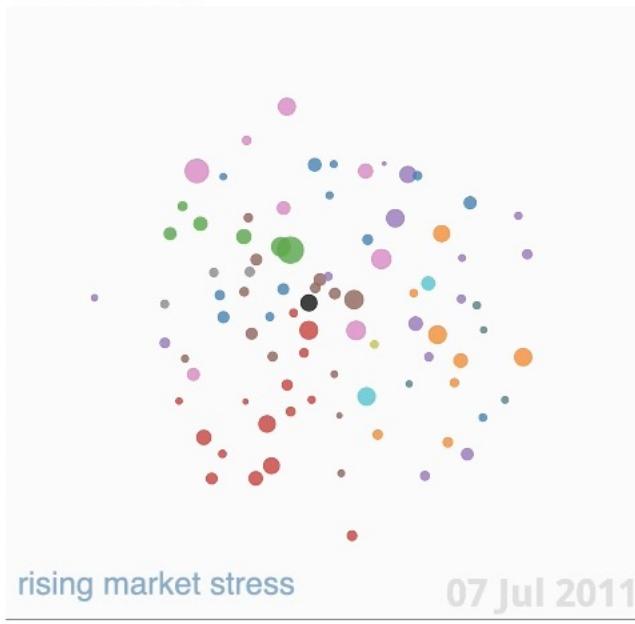
Hunting High and Low: Visualising Shifting Correlations in Financial Markets

Peter Simon^{1,2}, Cagatay Turkay¹

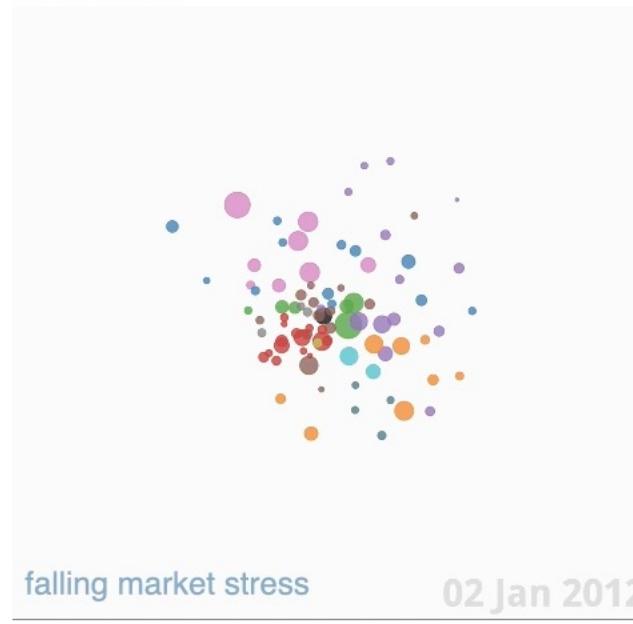
¹ City, University of London

² Scaridae Analytics

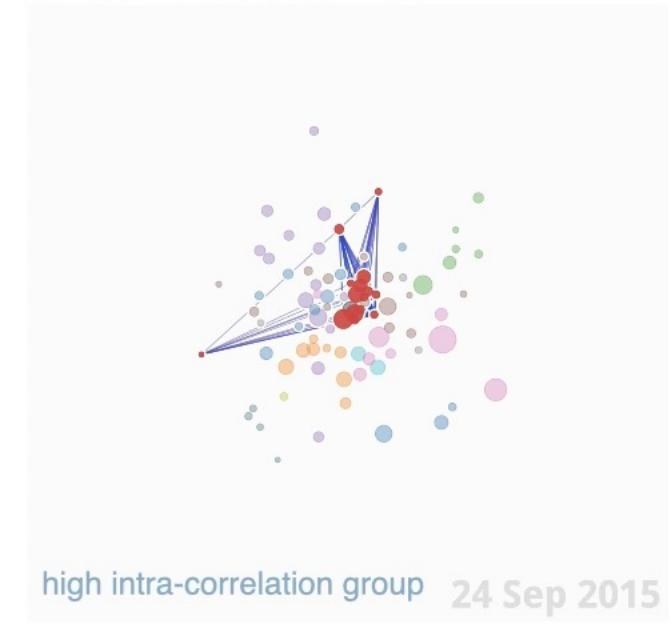
CLENCHING



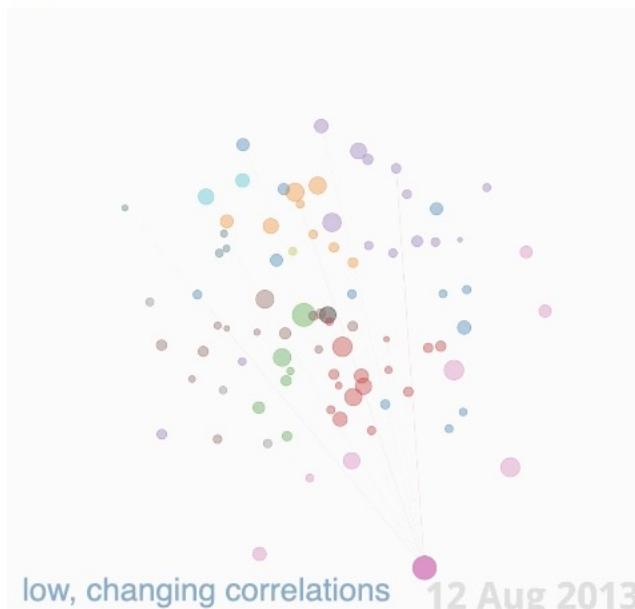
RELAXING



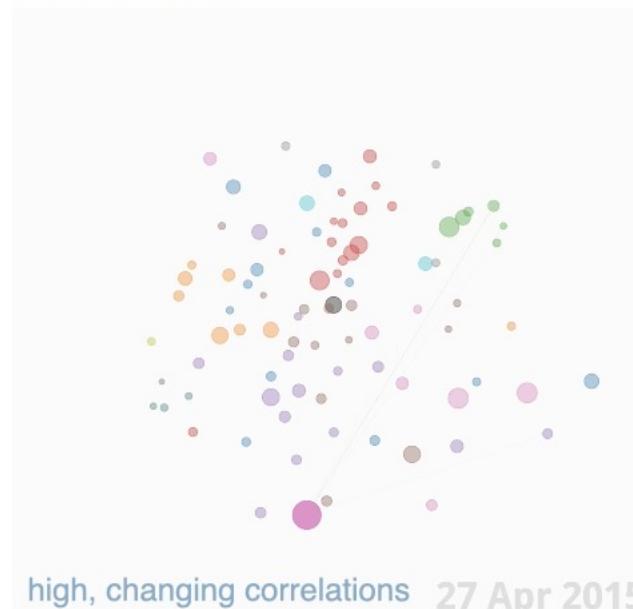
CLUSTERING/BUNCHING



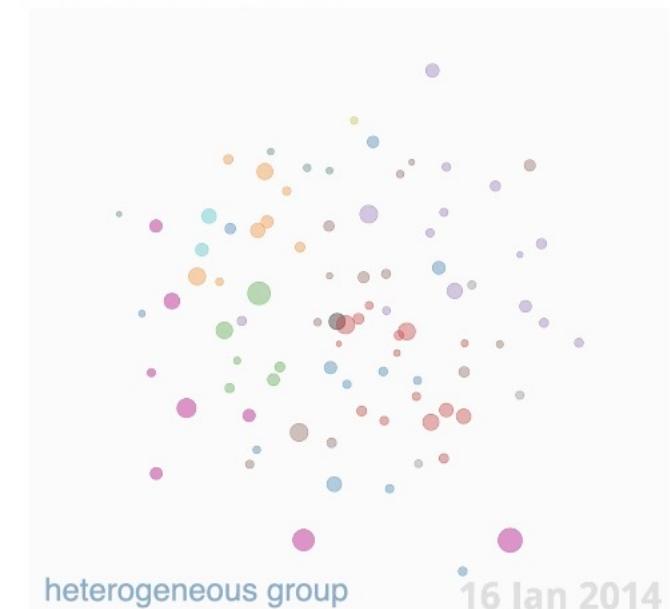
ORBITING



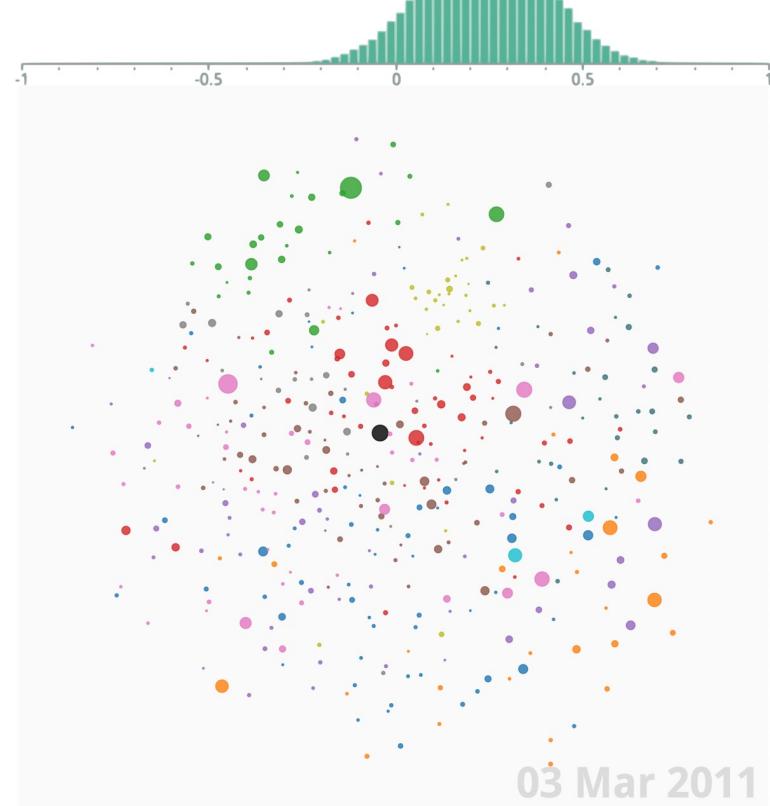
MEANDERING



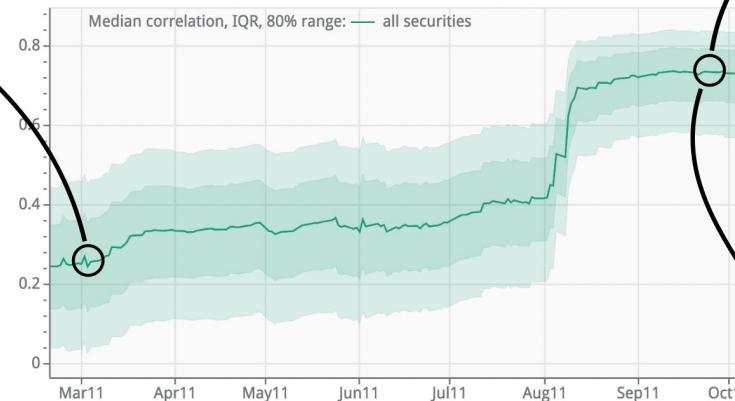
BROAD DISPERSION



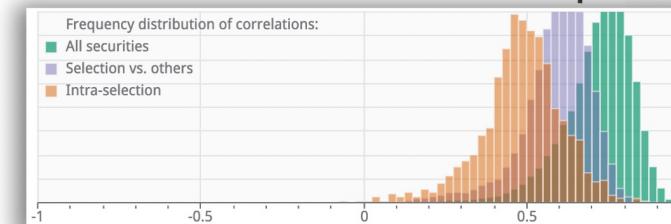
A. "Relaxed" No market stress



Intra-market correlations rose in 2011, jumping higher in August as concerns grew about banks' exposure to a possible Greek default and stock markets around the world wobbled



B. "Clenched" Severe market stress



Not all stocks moved in unison; the swarm's outliers were less correlated, both amongst themselves (orange) and with the rest of the market (blue)

- Index
- Financials
- Telecommunication Services
- Industrials
- Energy
- Consumer Staples
- Consumer Discretionary
- Real Estate
- Materials
- Information Technology
- Health Care
- Utilities

P3: generate **new perspectives** to “see” data and models

A few ideas to remember

- “Spaces” help us think, simplify, relate and also visualise Algorithms to provide alternatives, visualisation to link and relate them
- Always think about semantics – what does this “space” mean, what are the axes?

Four key practices in Visual Data Science

... facilitated by interaction and visualisation

P1: interrogate, relate, compare **variation, co-variation, deviation**

P2: interact with data and computational tools

P3: generate **new perspectives** to “see” data and models

P4: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

Phenomena is multi-faceted .. you need many views

their stories of what happened are very different.

What I call the Rashomon Effect is that there is often a multitude of different descriptions [equations $f(\mathbf{x})$] in a class of functions giving about the same minimum error rate. The most easily understood example is subset selection in linear regression. Suppose there are 30 variables and we want to find the best five variable linear regressions. There are about 140,000 five-variable subsets in competition. Usually we pick the one with the lowest residual sum-of-squares (RSS), or, if there is a test set, the lowest test error. But there may be (and generally are) many five-variable equations that have RSS within 1.0% of the lowest RSS (see Breiman, 1996a). The same is true if test set error is being measured.

So here are three possible pictures with RSS or test set error within 1.0% of each other:

Picture 1

$$y = 2.1 + 3.8x_3 - 0.6x_8 + 83.2x_{12} - 2.1x_{17} + 3.2x_{27},$$

Picture 2

$$y = -8.9 + 4.6x_5 + 0.01x_6 + 12.0x_{15} + 17.5x_{21} + 0.2x_{22},$$

Picture 3

$$y = -76.7 + 9.3x_2 + 22.0x_7 - 13.2x_8 + 3.4x_{11} + 7.2x_{28}.$$

Which one is better? The problem is that each one tells a different story about which variables are important.

The Rashomon Effect



from: <https://www.bfi.org.uk>

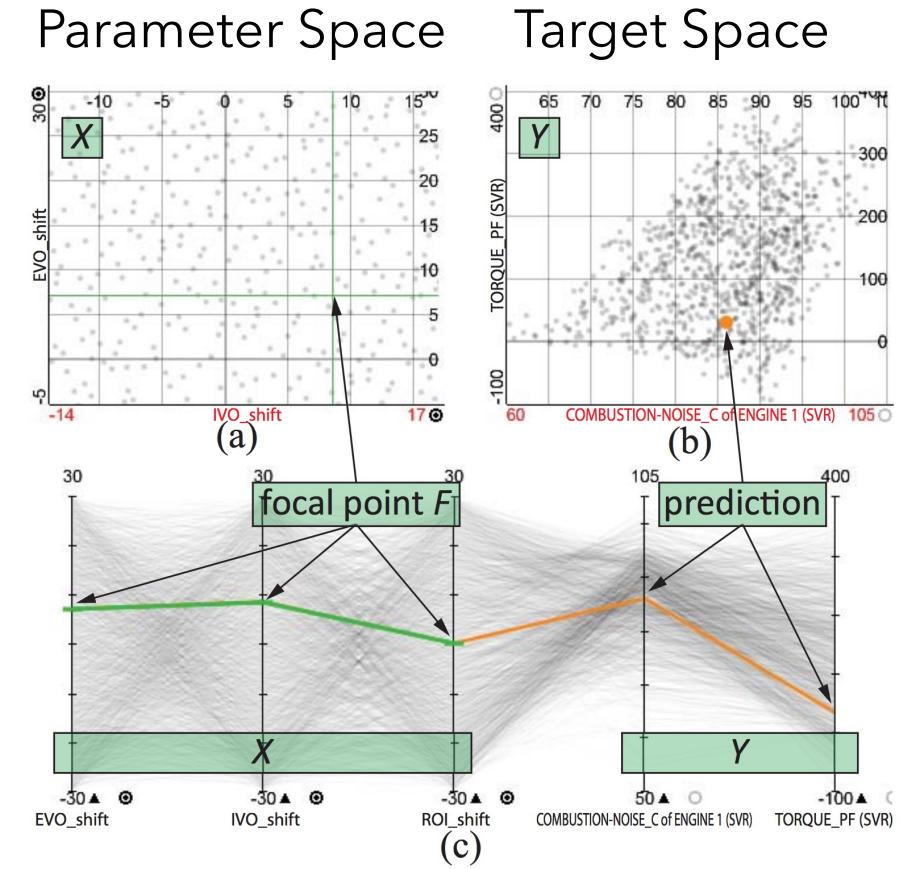
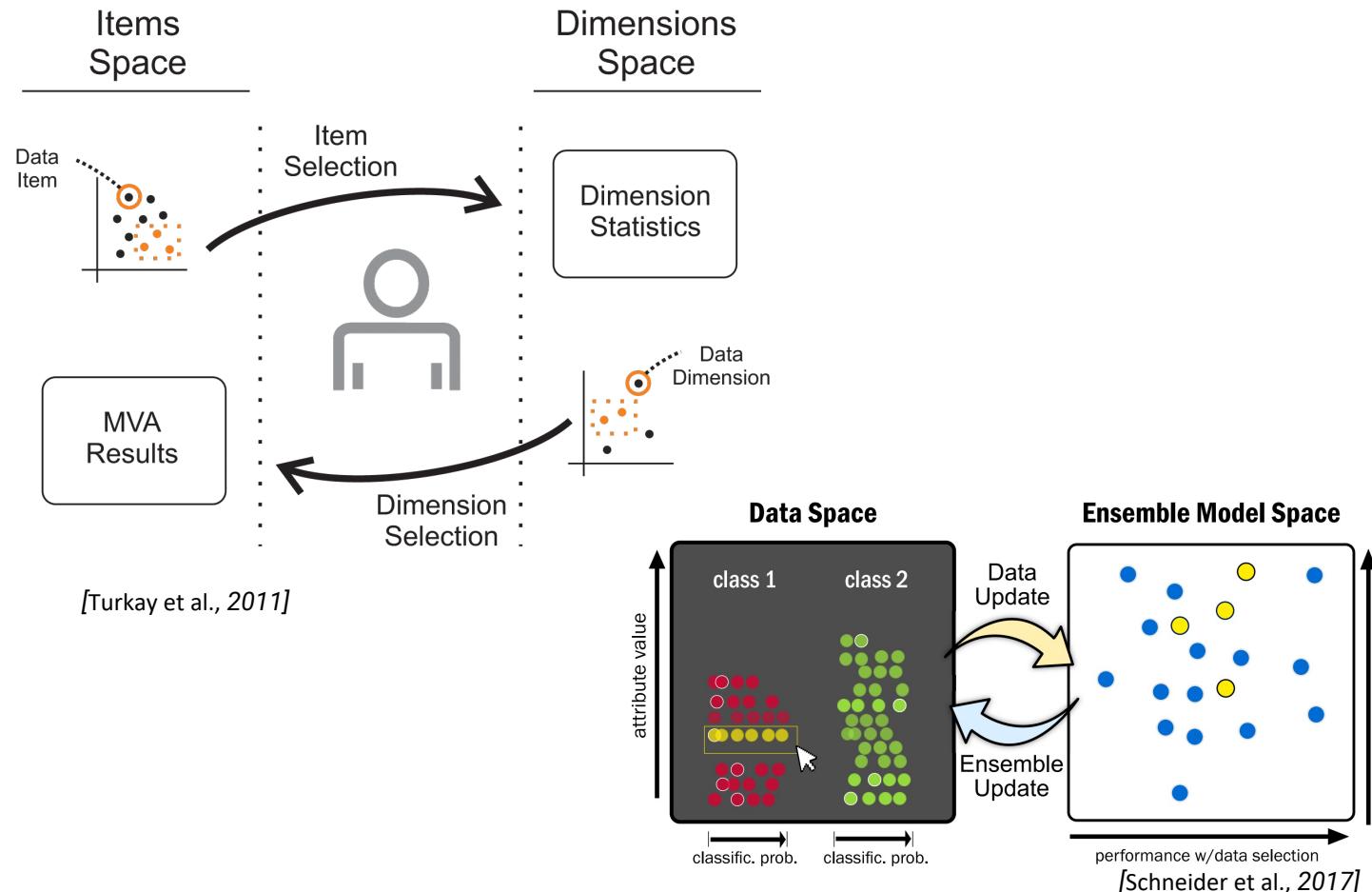
Statistical Modeling: The Two Cultures

Leo Breiman

Statistical Science
2001, Vol. 16, No. 3, 199–231

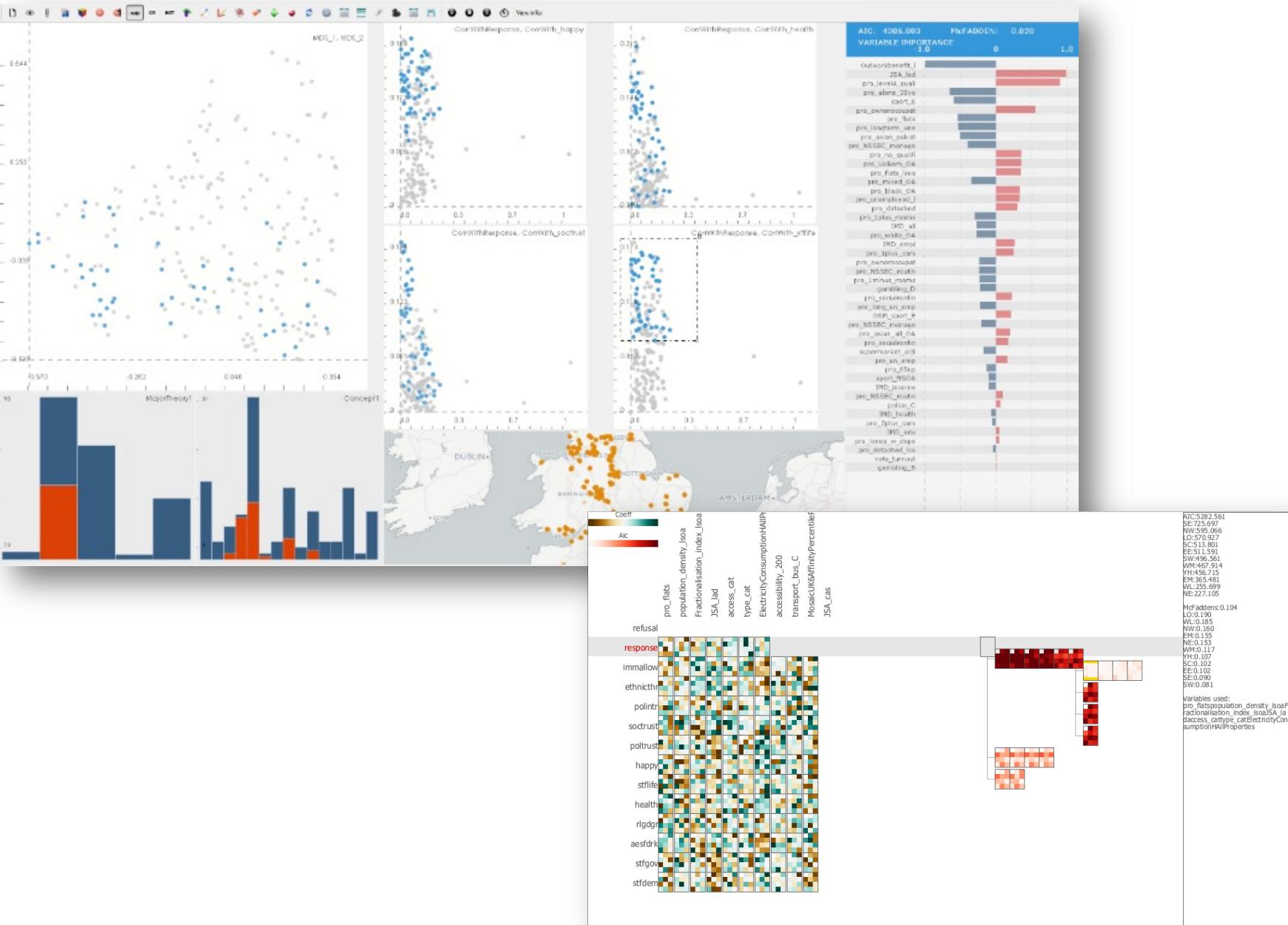
Multiple (data, model, ...) spaces

Building new semantic links between data/model/derived spaces



Enhancing a Social Science Model-building Workflow with Interactive Visualisation

Turkay, C., Slingsby, A., Lahtinen, K.,
Butt, S., & Dykes, J., ESANN 2016 (&
Neurocomputing 2017)

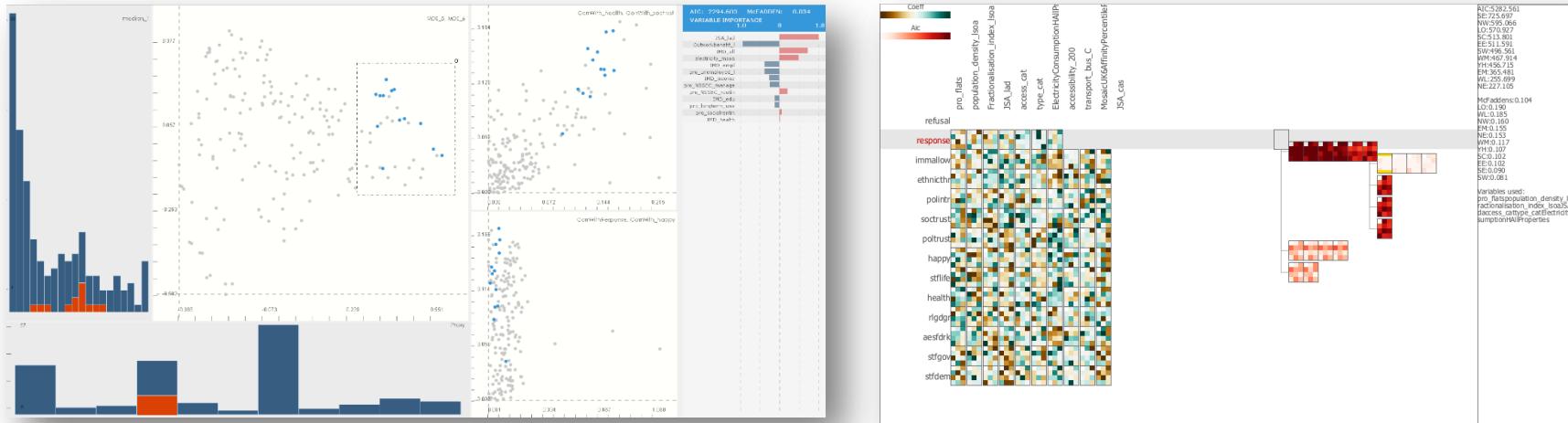


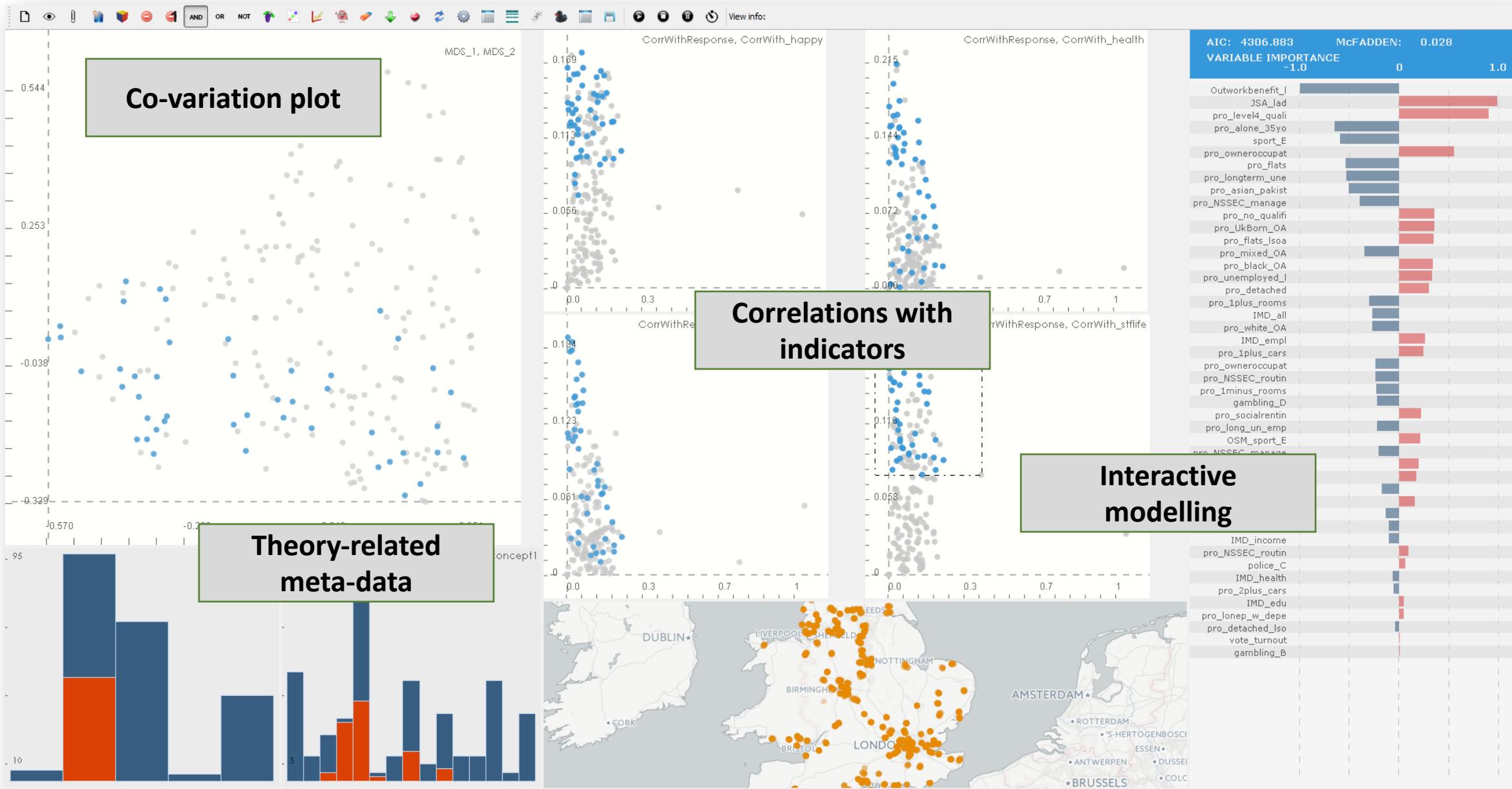
AddResponse - Details

- European Social Survey (ESS) UK 2012 - 13
- 4,520 households
- linked to auxiliary data from:
 - administrative sources
 - commercial consumer profiling
 - open-source data
- 401 auxiliary variables
- 32 survey response variables
(only for the respondents)



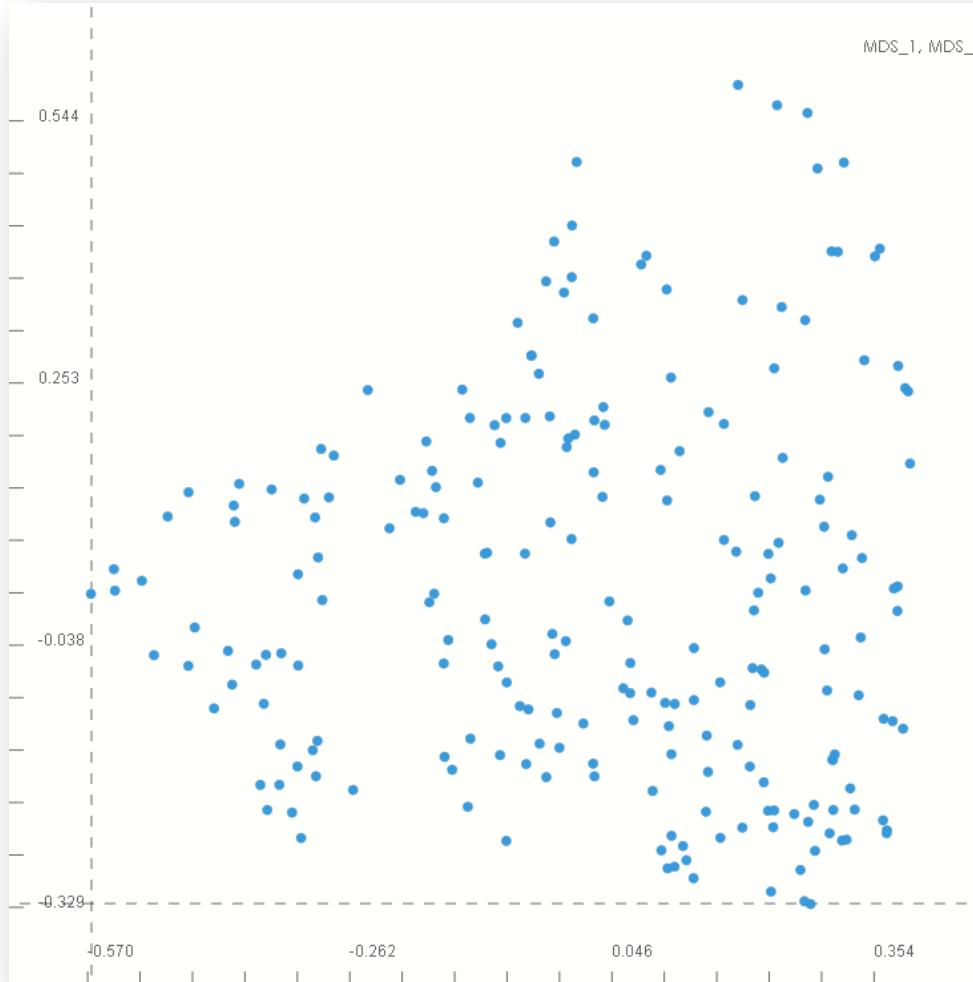
The aim is to investigate the relation of **auxiliary data** (i.e., 401 variables) to **survey non-response** and how these relate to the **indicators in survey responses** (i.e., 32 variables).







Exploring variables: Covariation & Correlation Maps

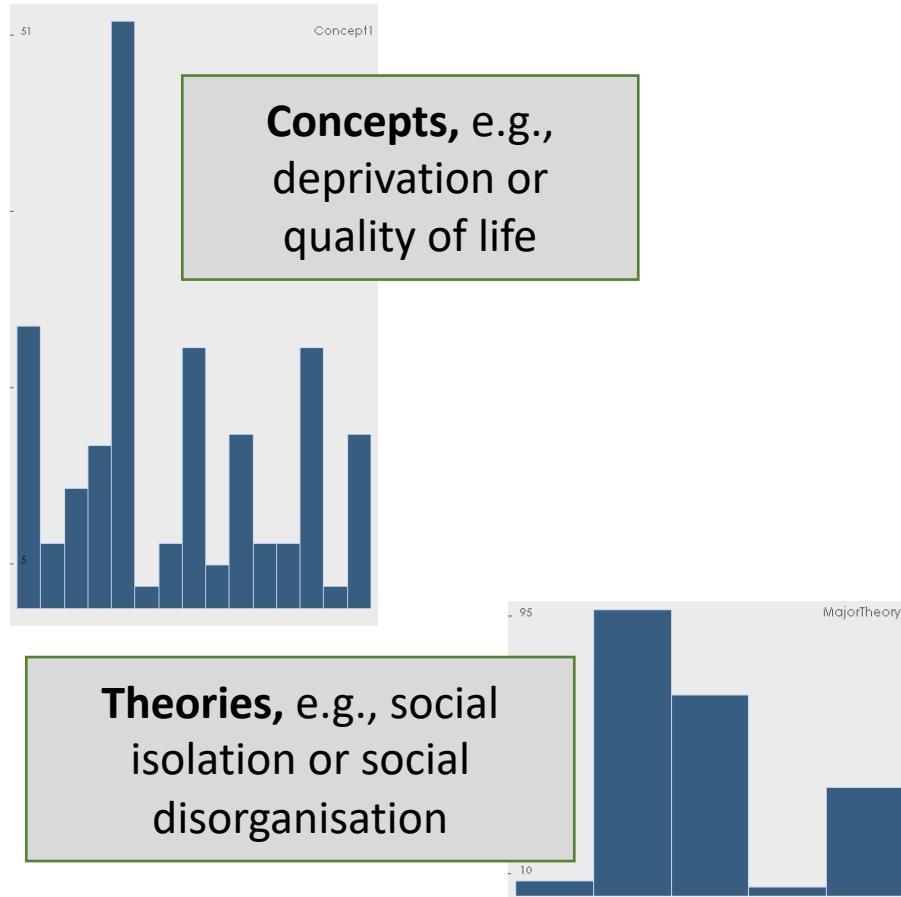


Compute **pairwise correlation**
within all 401 variables

Use this as a **distance matrix** and
project to 2D (using MDS)

Visualise on a scatterplot where
each point is a variable

Incorporating Theory-related data



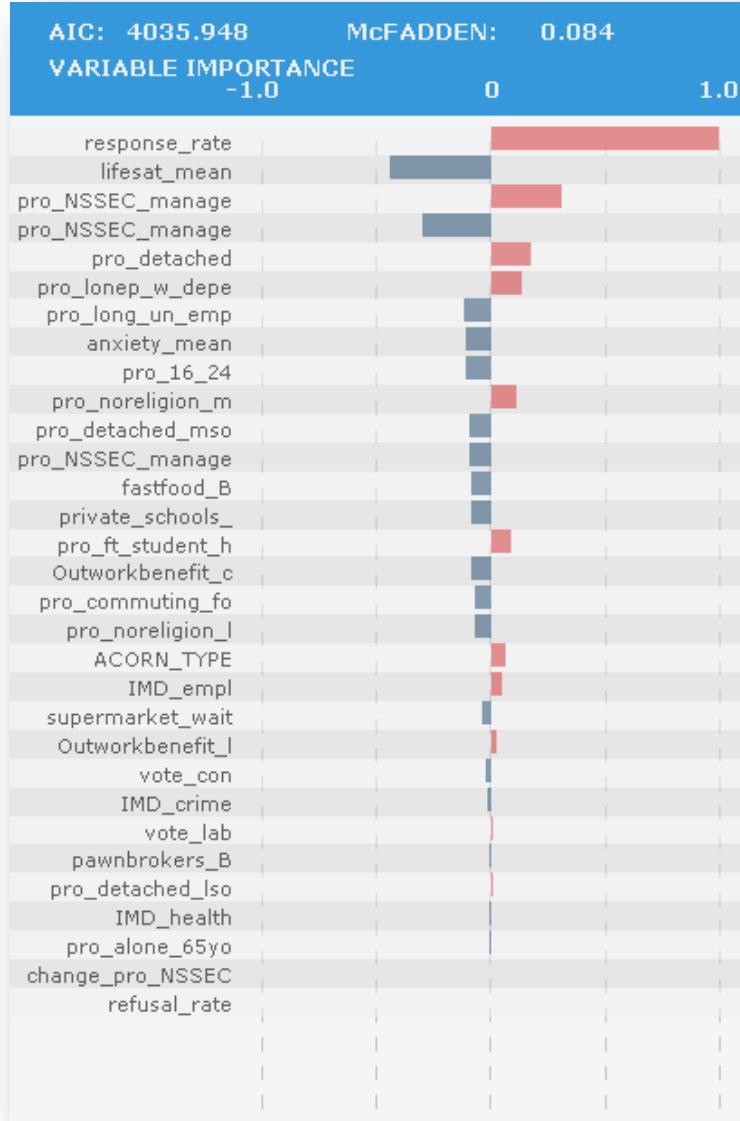
Associate variables to social-science
concepts and **theory**

Concepts relate to theories

Variables act as **proxies for concepts**

Use these as meta-data on variables
and visualise through histograms

Interactive model building

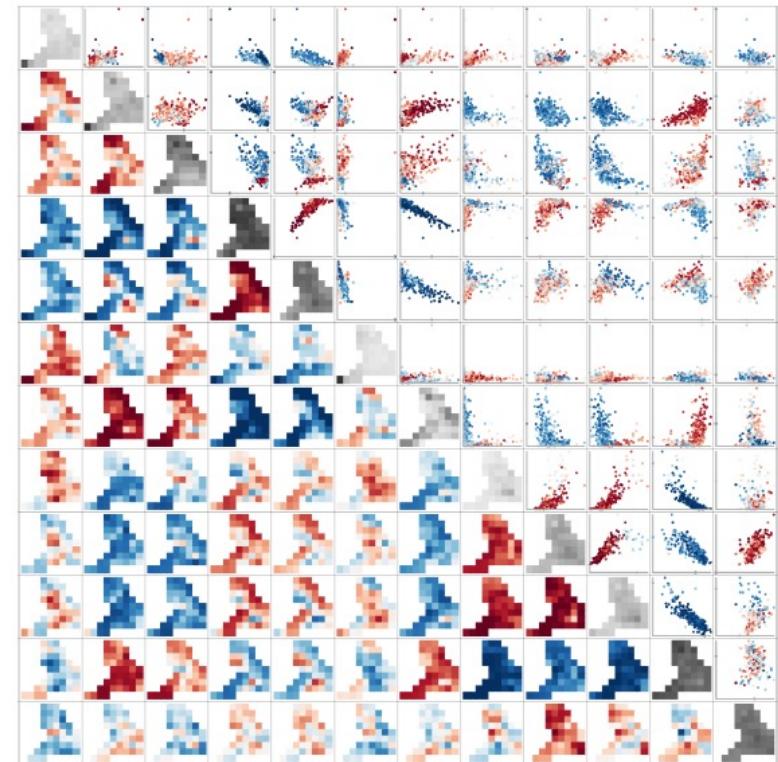
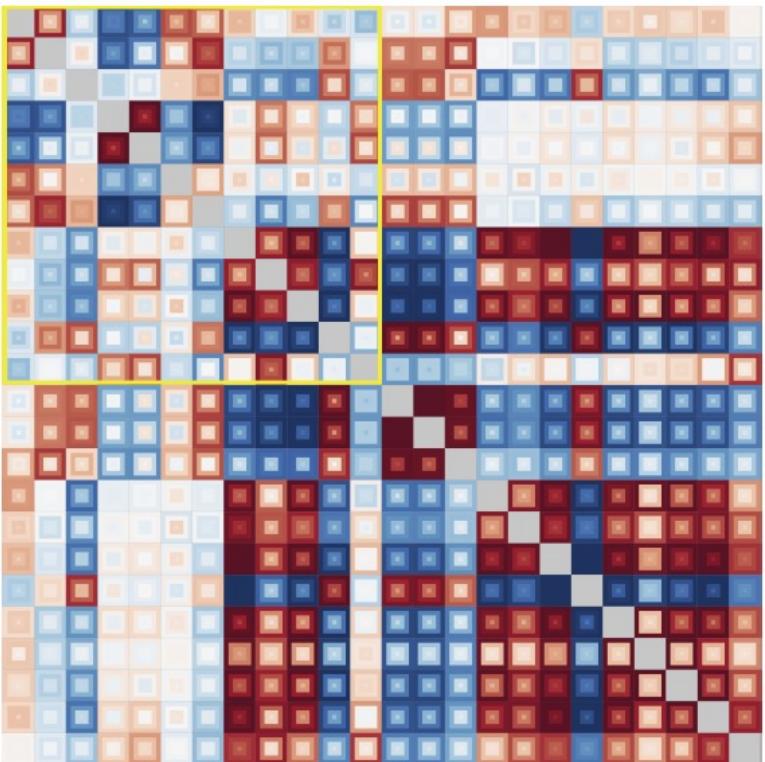
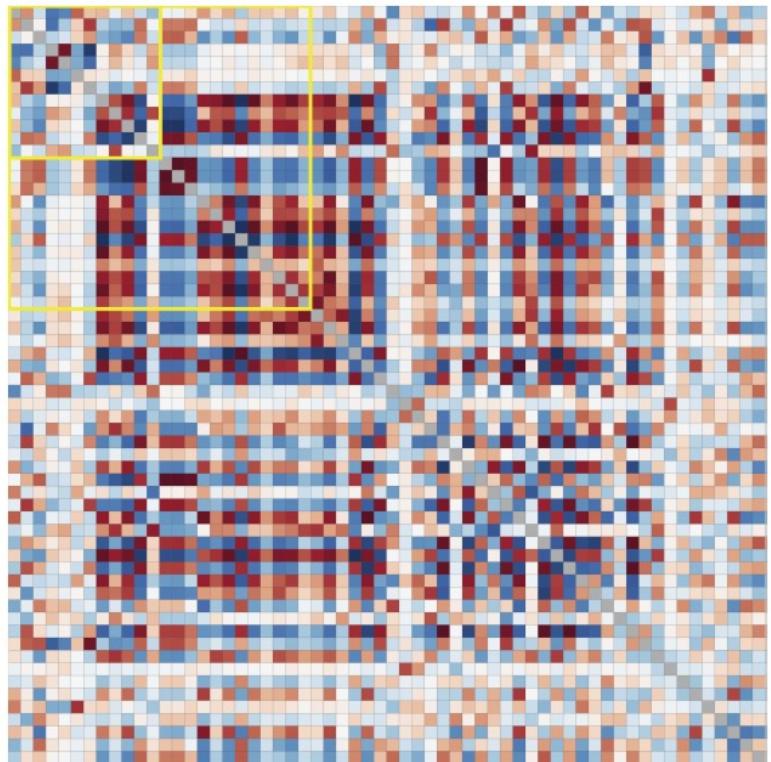


Embedded Logistic Regression

Variable weights visualised
“live”

Model fit scores indicated –
immediate feedback on
model quality

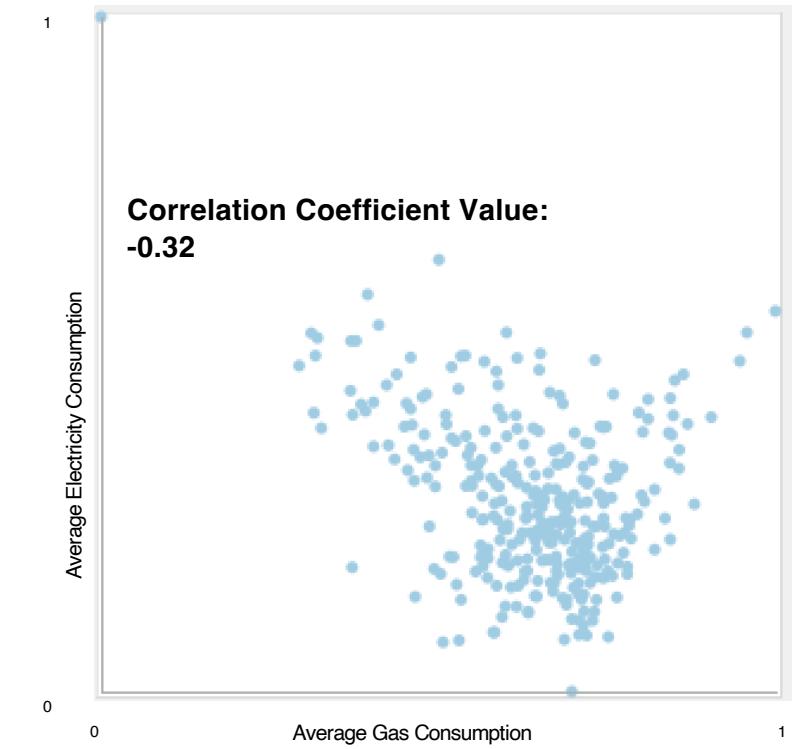
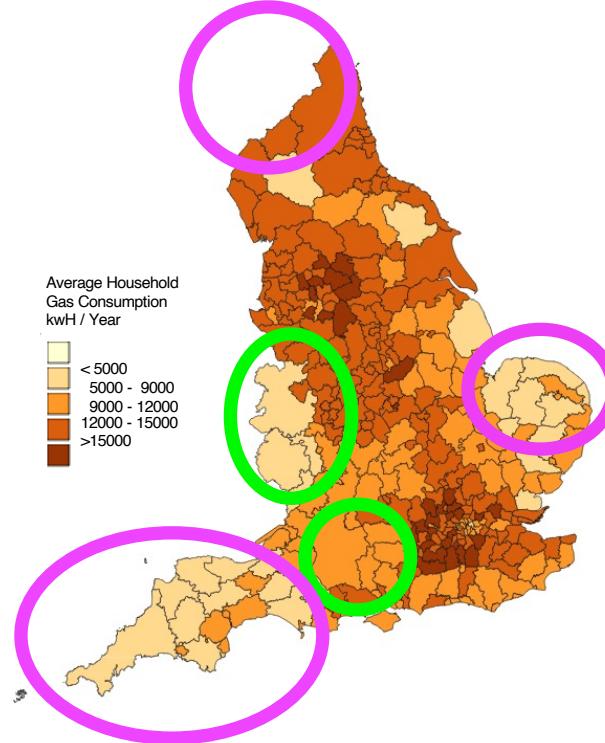
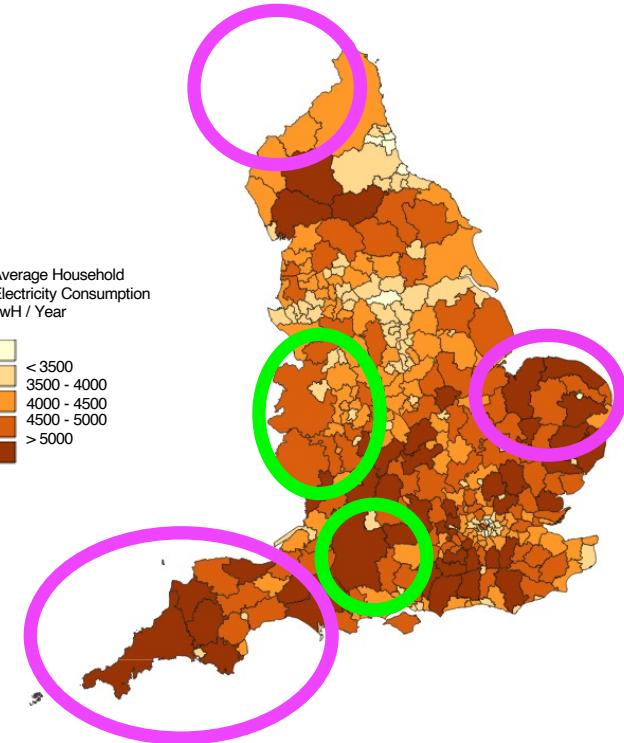




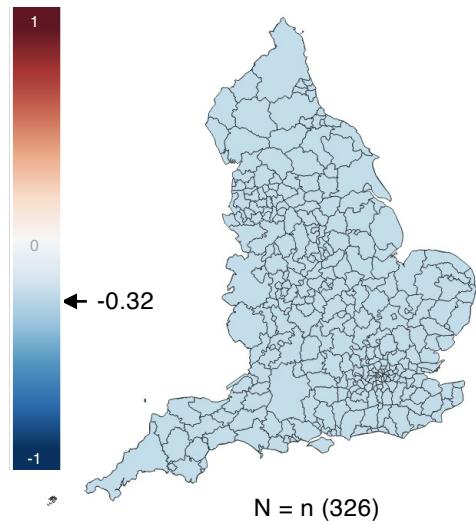
Visualizing Multiple Variables Across **Scale** and **Geography**

Goodwin, S., Dykes, J., Slingsby, A. & Turkay, C. (2016). Visualizing Multiple Variables Across Scale and Geography. InfoVis 2015.

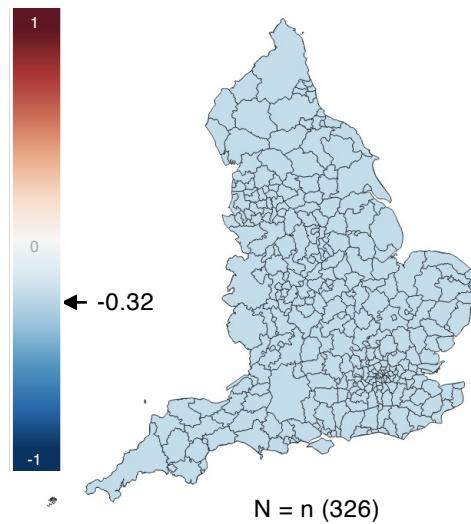
Bivariate Comparison



Local Correlation



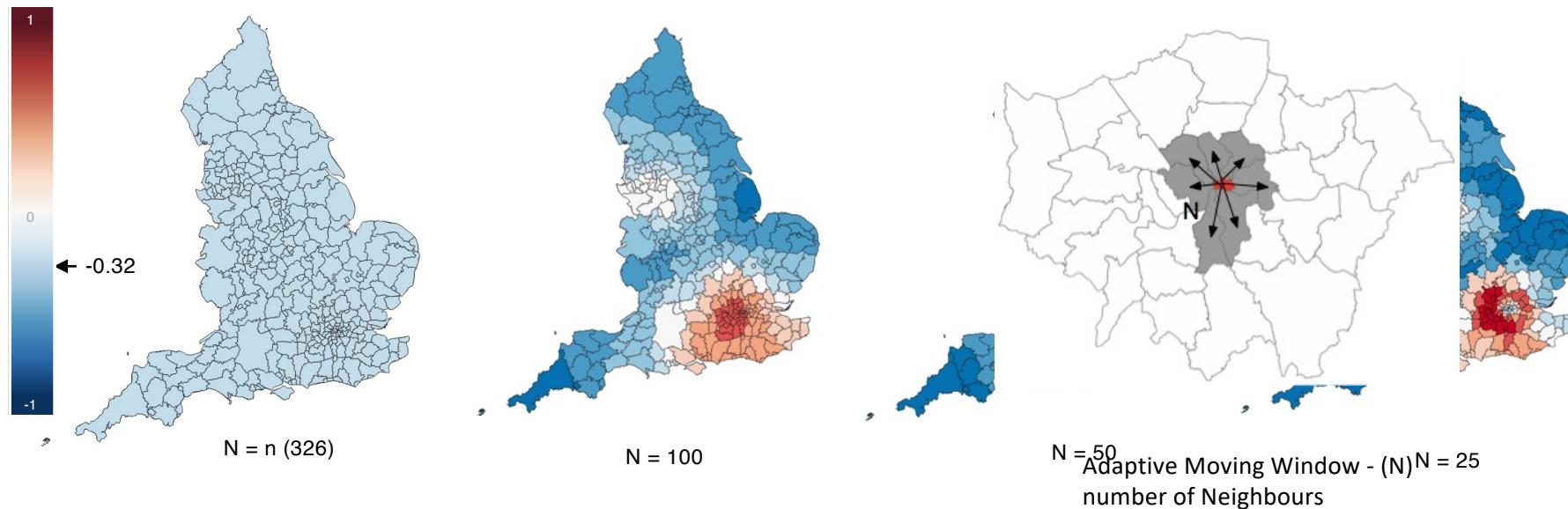
Local Correlation



Adaptive Moving Window - (N)
number of Neighbours

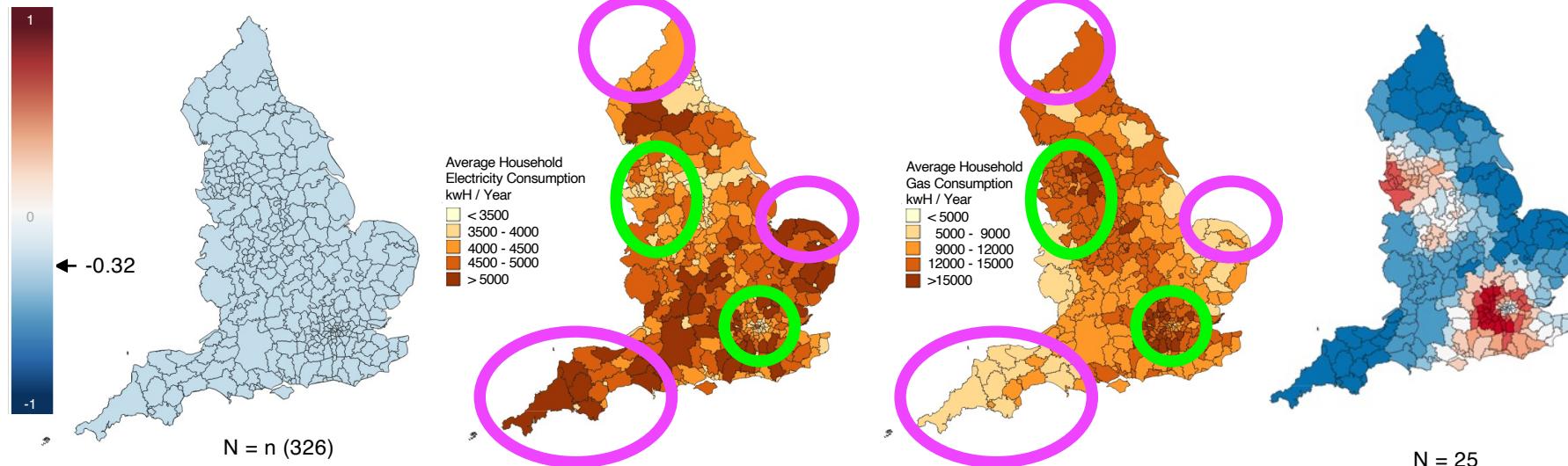
Brunsdon, C. Fotheringam, A. S. and Charlton, M. E. (1996) Geographically Weighted Regression: a method for exploring spatial nonstationarity, *Geographical Analysis*, 28(4): 281-298

Local Correlation

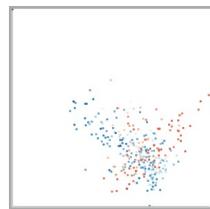


Brunsdon, C. Fotheringham, A. S. and Charlton, M. E. (1996) Geographically Weighted Regression: a method for exploring spatial nonstationarity, *Geographical Analysis*, 28(4): 281-298

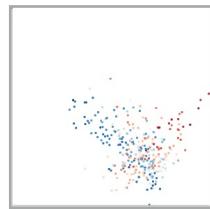
Local Correlation



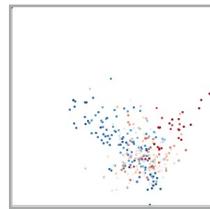
N = 25



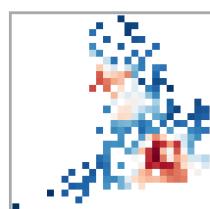
N = 50



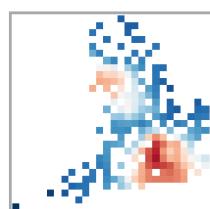
N = 100



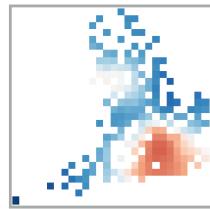
N = 25



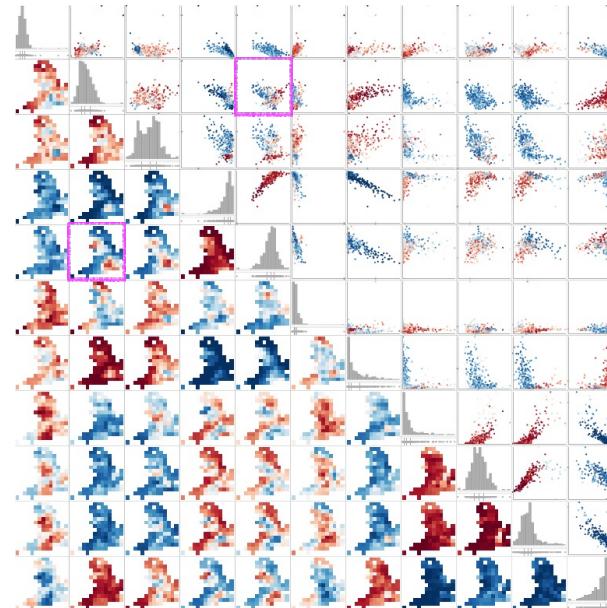
N = 50



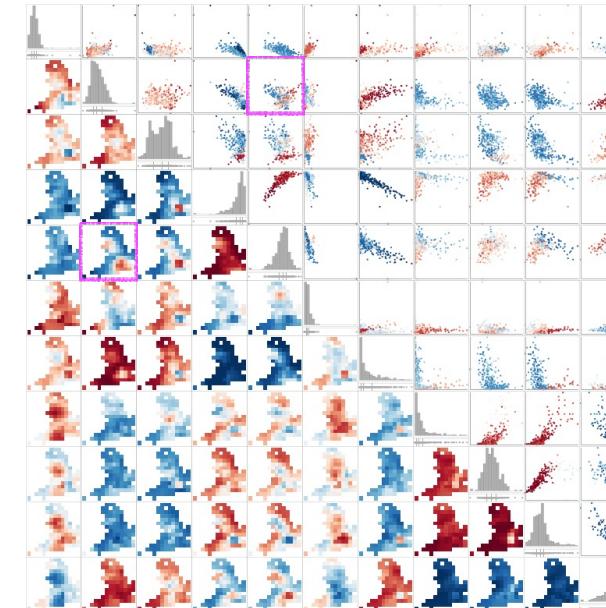
N = 100



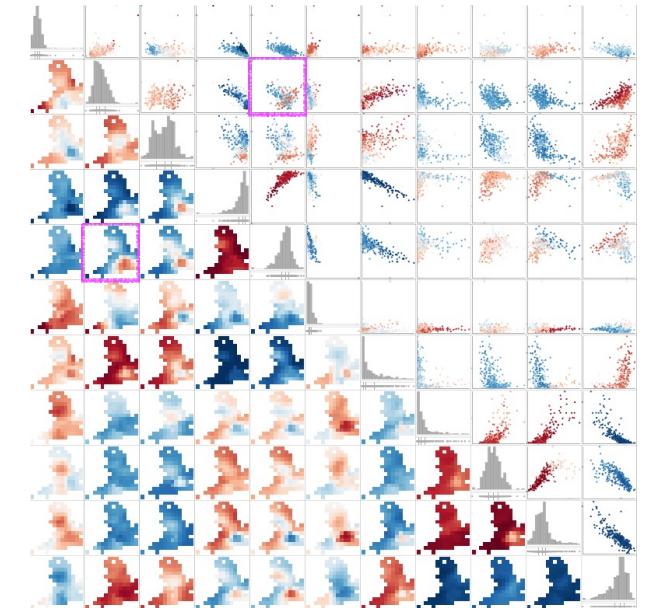
N = 25



N = 50

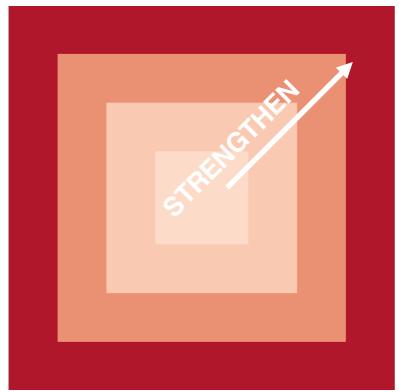
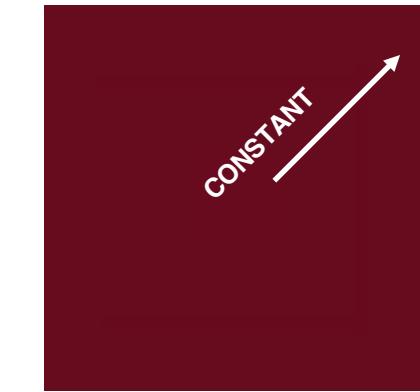
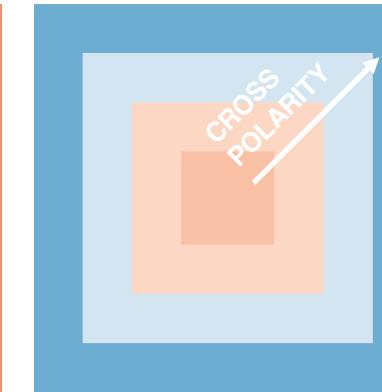
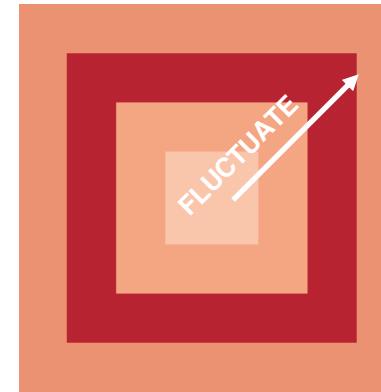


N = 100

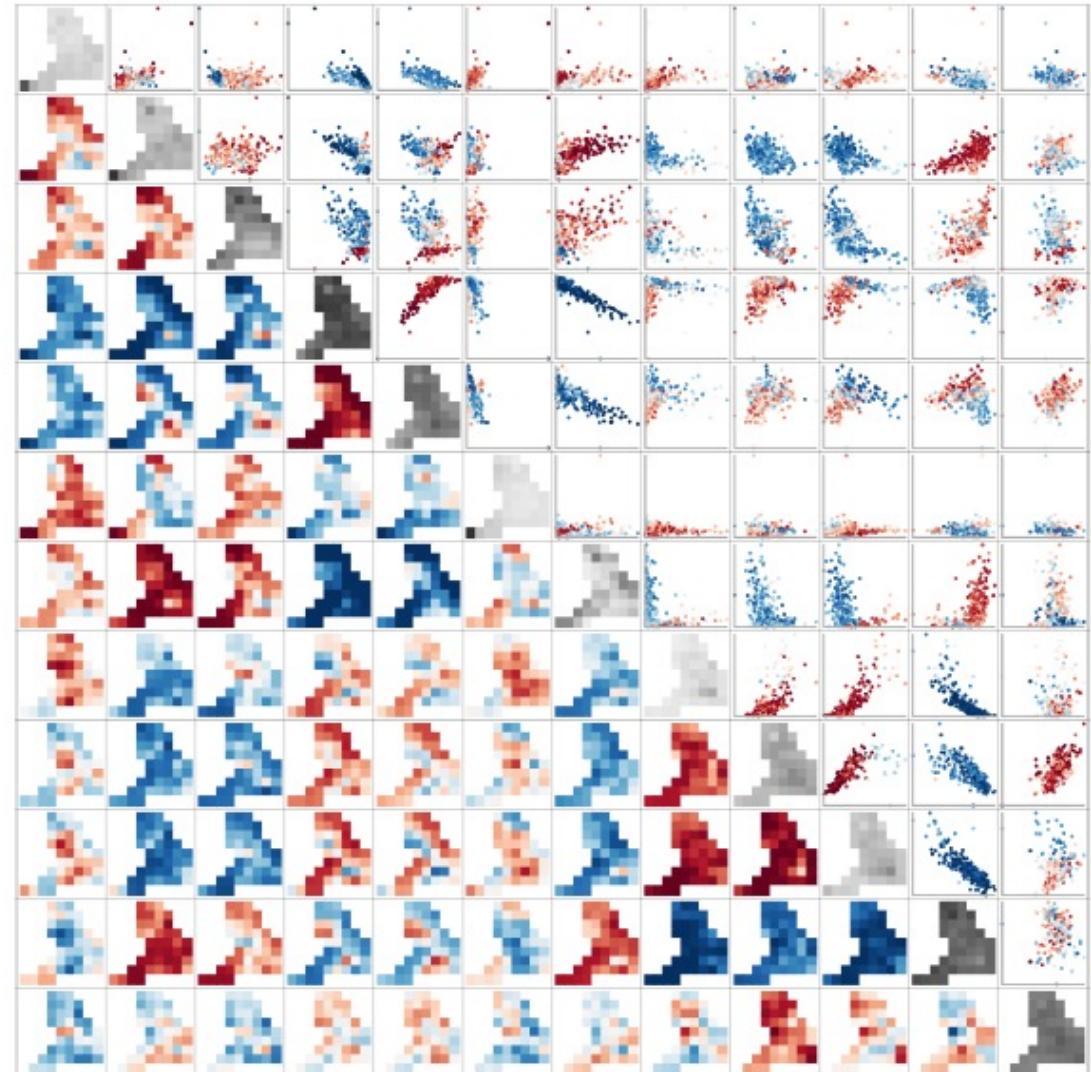
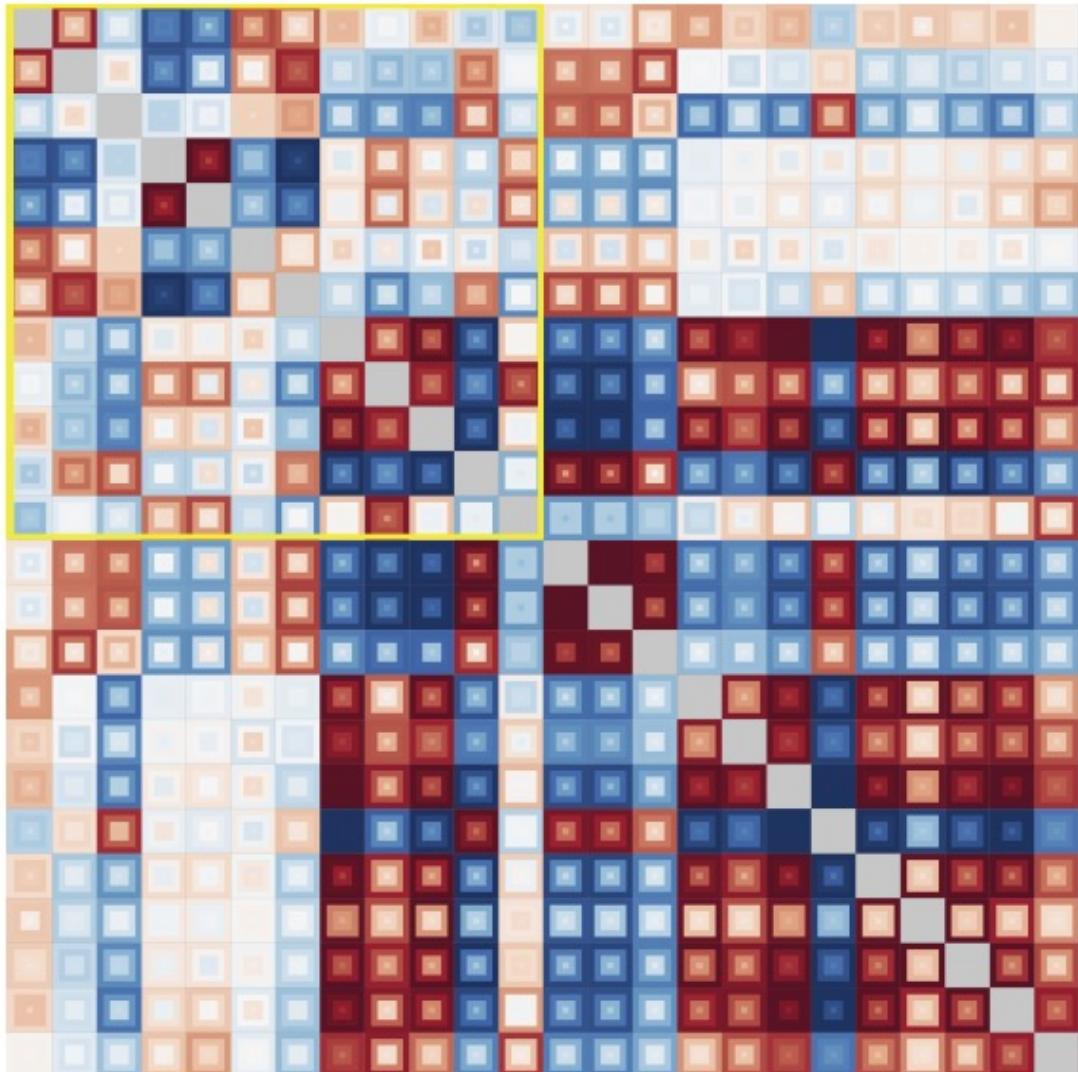


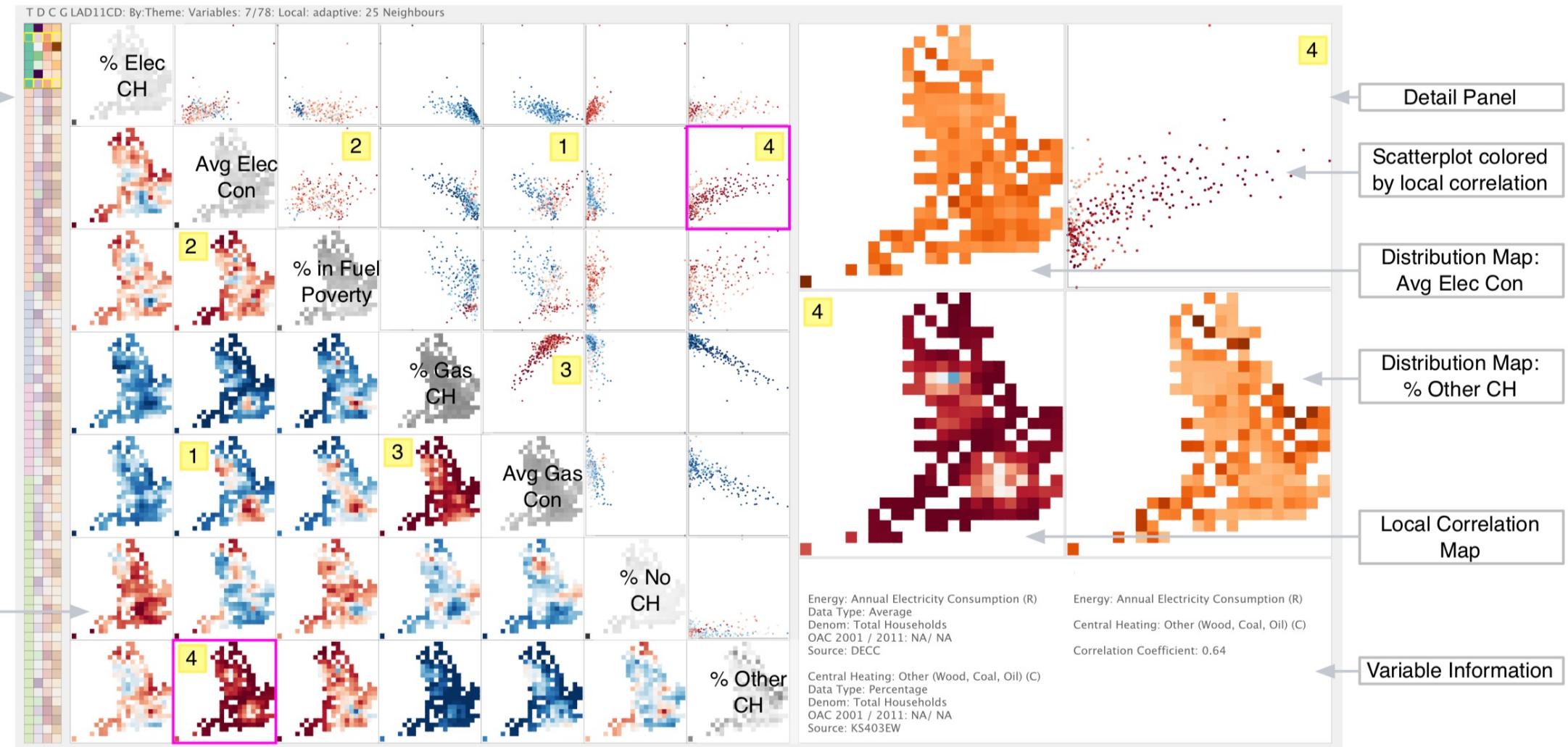
Many variables?

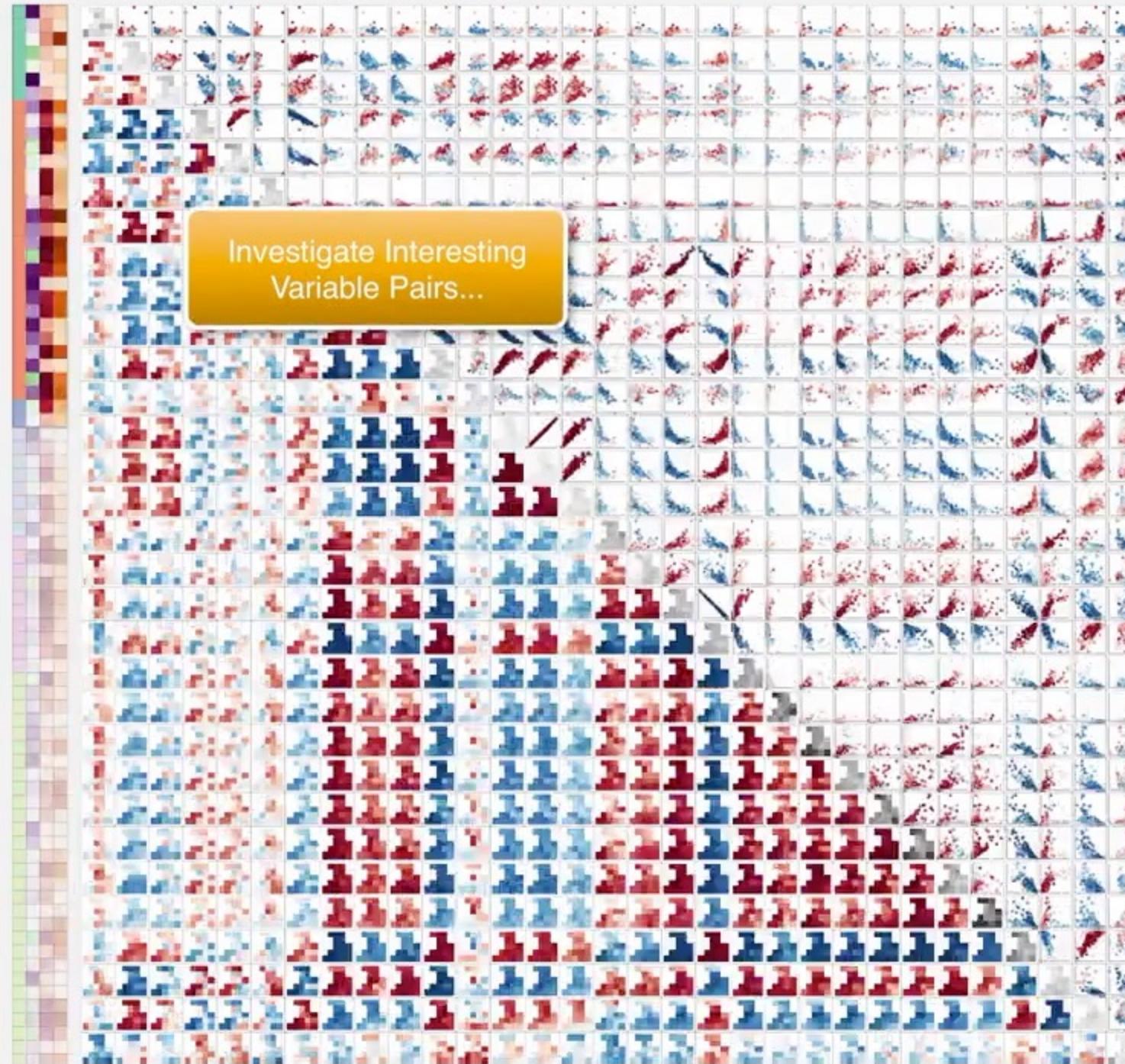
Multiple Scales Concurrently – Sensitivity



Many variables, many scales







P4: analyse **several aspects concurrently** (e.g., datasets, scales, parameters, algorithms ...)

A few ideas to remember

- Relate multiple-scales , multiple-variables, multiple datasets to infer new links, question assumptions, quality, (un)certainty
- Visualisation to concurrently display variations
- Interaction to “link” multiple facets/modalities might need some thinking

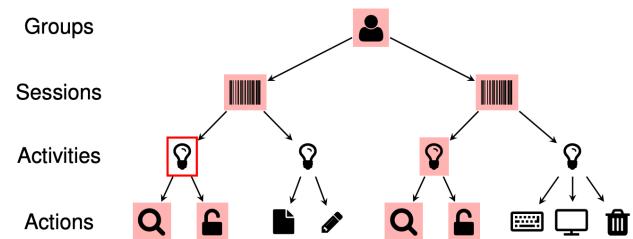


Fig. 7. An example of multi-semantics linking across four semantic levels. When an activity is mouse hovered (red border), all the same activities, sessions and users containing them, and individual actions within the activity are highlighted (pink backgrounds).

VDS: A FEW CLOSING REMARKS

Visualisation empowers you to ask & (hopefully) answer

HOW, WHY, WHAT IF

questions in data science practices ...

- How do these two models relate?
- Why is this a cluster?
- What if I leave this feature out?

....

.. through these questions ...

VISUALISATION
TO “**BETTER**”
DATA SCIENCE

RICH,
ENGAGED,
REFLECTIVE,
INFORMED

*dialogue with
data & models*

Next up .. VDS DIY

VDS DIY BLOCK #1

Thinking Visually

Part-1: Re/De-constructing Visualisations

- Introduction to visualisation basics
- Moving on to the visualisation grammar
- A "language" of visualisation
- Altair -- a versatile visualisation library in Python

Part-2: Foundational VDS Techniques

- Visualisation for interrogating the data
 - Conformity, Outliers, Shapes
- Multiple perspectives
 - Using small multiples

VDS DIY BLOCK #2

Visual Data Science -- Getting "involved"

Part-1: Working with models

- Working with models visually
 - Interpreting clustering
 - Working with projections for high-dimensional data
- Beyond accuracy

Part-2: Making it interactive

- Bringing in interaction
- Linking it up
- Going back to the data

VDS DIY BLOCK #3

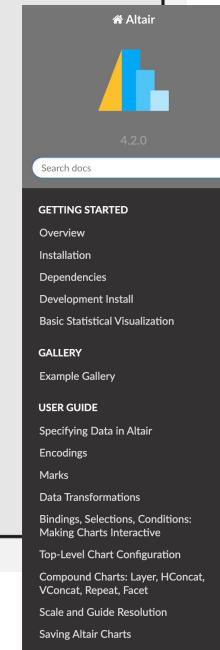
Doing (V)DS responsibly

Part-1: Being aware

- Dance of the p-values
- Maps and colouring

Part-2: Communicating openly

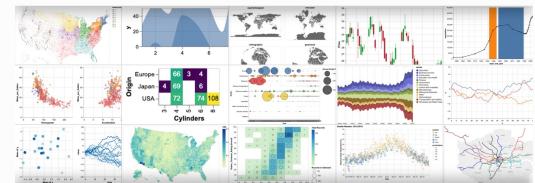
- Thinking about the narrative
- Communicating responsibly



» Altair: Declarative Visualization in Python

[View page source](#)

Altair: Declarative Visualization in Python



Altair is a declarative statistical visualization library for Python, based on [Vega](#) and [Vega-Lite](#), and the source is available on [GitHub](#).

With Altair, you can spend more time understanding your data and its meaning. Altair's API is simple, friendly and consistent and built on top of the powerful [Vega-Lite](#) visualization grammar. This elegant simplicity produces beautiful and effective visualizations with a minimal amount of code.

Getting Started

- Overview
- Installation
- Dependencies
- Development Install
- Basic Statistical Visualization

Gallery

- Example Gallery

- Hands-on exercises to explore some of these practices/ideas
- Will use Altair in Jupyter Notebooks to visualise and bring interactivity