# Modelling clusters from the ground up: a web data approach

Christoph Stich, Emmanouil Tranos & Max Nathan

University of Bristol, Alan Turing Institute
**e.tranos@bristol.ac.uk**, @EmmanouilTranos

# Introduction

## A lot of theoretical and empirical work on clusters

- Urban econ & econ geog (micro-foundations, MAR vs. Jacobs)
- Evolutionary perspectives (path dependency)
- Globalisation scholars (global value chains / production networks)
- Temporary / online collaboration tools

## Some basic questions still unresolved

- E.g. MAR vs. Jacobs; feasibility of cluster policy; appropriate policy mix
- Hard-to-fix empirical challenges:
- Data / economic activity scale mismatch (MAUP)
- SIC lag behind real-world industrial evolution
- Defining clusters based on industries instead of activities (e.g. fintech or cleantech)
- Tradeoffs between richness and reach of data

## Contribution

- A new approach to analyse clusters from the bottom up
- Over time
- Web data and data science methods
- Empirical cluster research challenges (MAUP, SIC, richness/reach tradeoff)
- Shoreditch: East London Tech City aka Silicon Roundabout

# Empirical strategy

# Web data

- *Archived*, commercial websites 2000-2012
- Geolocated in Shoreditch, London
- Flexible approach in exploring economic activities and their dynamics
- Readily available, cheap to obtain and extensive in terms of the theme and population coverage
- Under-explored, public domain data

# Methods

- Data cleaning: create a subset of business websites in Shoreditch
- Spatial analysis for interesting outliers
- Topic modelling: Latent Dirichlet Allocation (LDA)
- Extract bundles of economic activities (topics)
- Extract the key terms of every topic
- Bottom up classification *vs*. top-down SIC

# Data

# Web data: The Internet Archive

- The largest archive of webpages in the world
- 273 billion webpages from over 361 million websites, 15 petabytes of storage (1996 -)
- A web crawler starts with a list of URLs (a seed list) to crawl and downloads a copy of their content
- Using the hyperlinks included in the crawled URLs, new URLs are identified and crawled (snowball sampling)
- Time-stamp

# Web data: The Internet Archive

# Web data: The Internet Archive

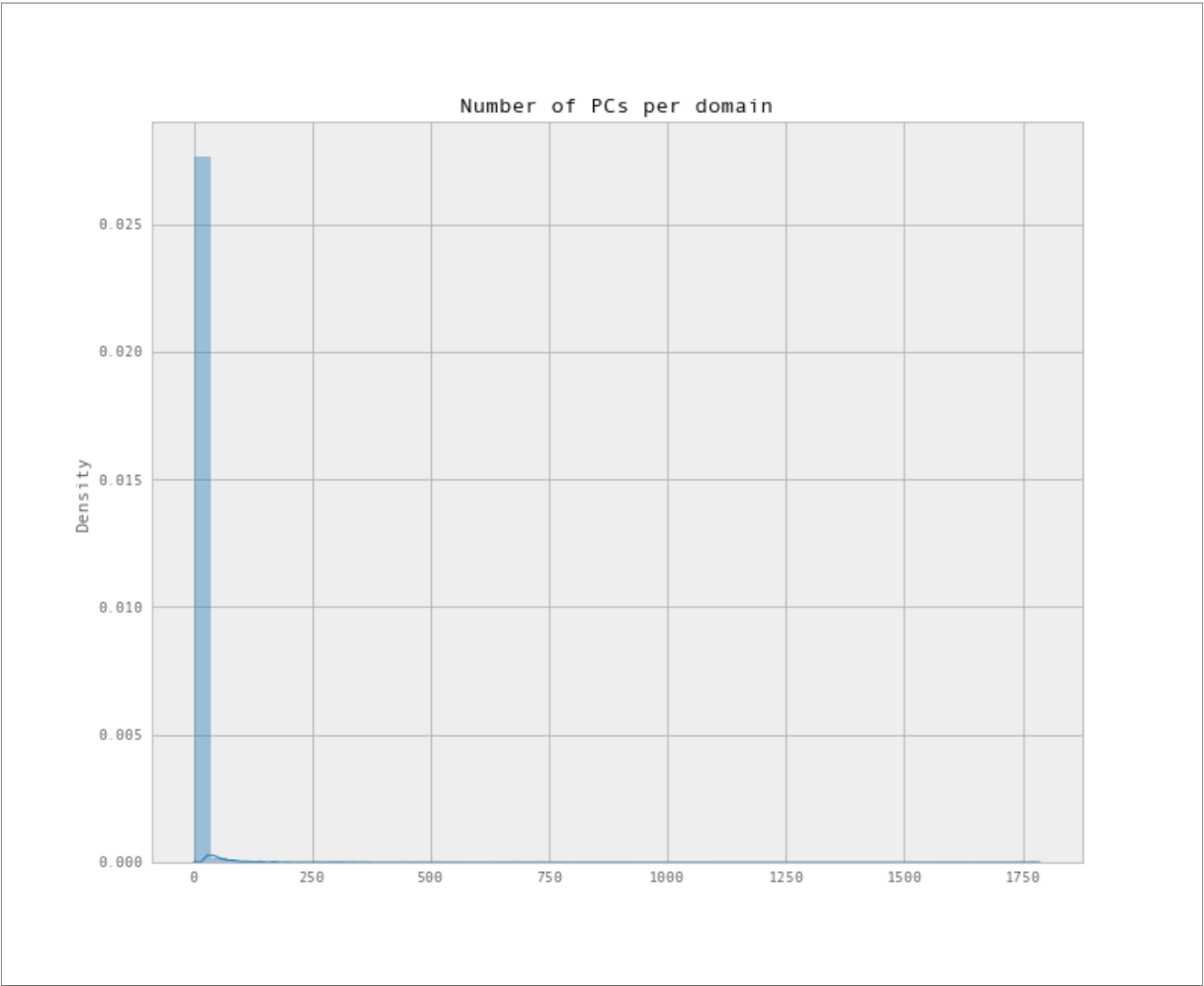# Web data: The Internet Archive

- JISC UK Web Domain Dataset: all archived webpages from the .uk domain 1996-2012
- Curated by the British Library
- Geoindex: a subset of the .uk archived webpages which contain a UK postcode
- circa 0.5 billion URLs with valid UK postcodes

20080509162138/**http://uk.eurogate.co.uk/contact_us** IG8 8HD

# Data cleaning

- All the archived .uk webpages
- Archived during 2000-2012
- Commercial webpages (.co.uk & .ltd.uk)
- A postcode *in the web text* within the Shoreditch area
- From webpages to websites: **http://www.website1.co.uk/webpage1** and **http://www.website1.co.uk/webpage2** are part of the **http://www.website1.co.uk**
- 1 *vs.* multuple postcodes in a website

# Data cleaning



Number of PCs per domain

# Data cleaning

- Right side: websites with a large number of postcodes (e.g. directories, real estate websites)
- Left side: websites with a unique postcode in Shoreditch

# Directory website with a lot of postcodes

# Website with a unique postcode in Shoreditch

# Data cleaning

- Current analysis: website with a *unique* postcode in Shoreditch
- 71% of all the archived, commercial, geolocated websites for 2010
- Sensitivity: repeat the analysis including websites with up to 9 postcodes, at least one within Shoreditch
- 95% of all the archived, commercial, geolocated websites for 2010
- We observe **economic activities** and **not firms** within industries
- Websites do not necessarily correspond to firm entities

# Results

# Spatial concentration



Distribution of observations

# Websites per postcode

# Outlier

# Digital squatting



**CAPITAL OFFICE**

HOME  OUR SERVICES  FAQ'S  BLOG  LOGIN

CALL US +44 (0) 207 566 3939

Airplane mode off

## Squatters – illegal use of our address

Help prevent fraud

### Illegal use of our mail box services

As one of the leading virtual office address providers in London, we are determined to prevent fraudsters and squatters who use our mailbox service for illegal purposes or without our consent. Squatting is a term whereby people use display our address on their correspondence materials such as websites, business cards, letters heads and have not been given any consent. In most circumstances they are using this for illegal activity.

**Latest Identified Squatters:**

4 Sold
Atrofi Design Ltd
Best Accessories UK
Case Stop Ltd
Centre for Medical Science
Control Your Credit UK
EHIC Services Europe
Fraser Tores Property
GCR Capital
Gettickets.co
Instant Lending
JPD Tree and Garden Services
Lloyd Loom Spalding
Prime Brokerz
Recruit Network
Status Hair
Zuum Hoverboards

ALI FINANCIAL SERVICES LTD

We are here!

Chat now

Selected Topic: [0]  [Previous Topic] [Next Topic] [Clear Topic]

Slide to adjust relevance metric:[2]

λ = 1

| 0.0 | 0.2 | 0.4 | 0.6 | 0.8 | 1.0 |

## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

1

3

2

4

Marginal topic distribtion

2%

5%

10%

## Top-30 Most Salient Terms[1]

| 0 | 1,000 | 2,000 | 3,000 | 4,000 |

design
job
web
properti
invest
websit
account
hotel
manag
recruit
market
cours
financi
busi
financ
club
holiday
art
shop
print
compani
employ
graphic
insur
brand
digit
train
health
site
car

Overall term frequency

Estimated term frequency within the selected topic

**1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)**
**2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)**

# Topics over time



Absolute prevalence topics

Legend:
- 0 design web websit servic art
- 1 hotel shop club holiday car
- 2 manag servic job busi properti
- 3 train cours health care support

# Topics over space

Selected Topic: 0 | Previous Topic | Next Topic | Clear Topic

Slide to adjust relevance metric:

λ = 1      0.0   0.2   0.4   0.6   0.8   1.0

## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

2

1

9

7

10  12
13
15

5   3

11

4

8

6

Marginal topic distribtion

2%

5%

10%

## Top-30 Most Salient Terms[1]

0    500   1,000  1,500  2,000  2,500  3,000

art
photograph
web
photographi
design
photo
artist
book
websit
light
music
print
card
site
model
galleri
search
engin
educ
agenc
host
anim
busi
exhibit
onlin
architectur
build
imag
internet
architect

Overall term frequency

Estimated term frequency within the selected topic

**1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)**
**2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)**

## Topics and keywords

- **Digital media**: design, brand, art, graphic, digit, print
  - Digital content creation; internet services and advertisement
  - publishing and performance arts; visual arts; photography services
  - interior design; garden services; home appliances
- **Visitor and leisure economy**: hotel, club, shop, holiday, car, travel
- **Financial and business service activities**: job, manang, properti, servic, busi, invest, account
- **Health and education**: cours, train, health, care, learn, test, treatment
- Topics: bundles of economic activities

# Traditional data

- Administrative data: UK registrar of companies
- SIC codes
- Plotting frequencies of SIC within Shoreditch
- Firms active 2000-2012

# Traditional data

| SIC Codes | Count | Description | Share |
|---|---|---|---|
| 70229 | 1134 | Management consultancy activities **other** than financial management | 0.201 |
| 64999 | 517 | Financial intermediation **not elsewhere classified** | 0.092 |
| 74909 | 387 | **Other** professional, scientific and technical activities **n.e.c.** | 0.069 |
| 68209 | 371 | **Other** letting and operating of own or leased real estate | 0.066 |
| 62012 | 326 | Business and domestic software development | 0.058 |
| 78109 | 185 | **Other** activities of employment placement agencies | 0.033 |
| 64209 | 171 | Activities of **other** holding companies **n.e.c.** | 0.030 |
| 56101 | 157 | Licensed restaurants | 0.028 |
| 59111 | 154 | Motion picture production activities | 0.027 |
| 69201 | 130 | Accounting and auditing activities | 0.023 |
| 71111 | 123 | Architectural activities | 0.022 |
| 43999 | 86 | **Other** specialised construction activities **n.e.c.** | 0.015 |
| 64205 | 85 | Activities of financial services holding companies | 0.015 |
| 93199 | 73 | **Other** sports activities | 0.013 |
| 56302 | 69 | Public houses and bars | 0.012 |
| 68201 | 67 | Renting and operating of Housing Association real estate | 0.012 |
| 69109 | 66 | Activities of patent and copyright agents; **other** legal activities **n.e.c.** | 0.012 |
| 59112 | 66 | Video production activities | 0.012 |

| SIC Codes | Count | Description | Share |
|---|---|---|---|
| 70221 | 65 | Financial management | 0.012 |
| 62011 | 64 | Ready-made interactive leisure and entertainment software development | 0.011 |
| 59113 | 63 | Television programme production activities | 0.011 |
| 71129 | 61 | **Other** engineering activities | 0.011 |
| 41201 | 58 | Construction of commercial buildings | 0.010 |
| 56102 | 56 | Unlicensed restaurants and cafes | 0.010 |
| 41202 | 47 | Construction of domestic buildings | 0.008 |
| 69202 | 45 | Bookkeeping activities | 0.008 |
| 64991 | 43 | Security dealing on own account | 0.008 |
| 58142 | 41 | Publishing of consumer and business journals and periodicals | 0.007 |
| 74209 | 40 | Photographic activities **not elsewhere classified** | 0.007 |
| 18129 | 40 | Printing **n.e.c.** | 0.007 |
| **Total** | | | **0.849** |

## Conclusions

- Modelling clusters and their dynamics *is not* a trivial problem
- Hard-to-solve empirical challenges
- Powerful and flexible approach
  - empirical challenges
  - implement key theoretical concepts (within-cluster co-location patterns, local distinctiveness, related / unrelated variety of activity, and cluster evolution)
- More informative than next-best analysis using open administrative data
- Detect unknown or emerging cluster formations