ATS1208 Digital Humanities: Concepts, tools and debates

Week 10: Impacts - Scholarship



dh@monash - ATS1208

Overview

- Space and digital scholarship
 - GIS
- Space as a variable
 - Comparing two places: Thomas and Ayers
 - Comparing many places: Goncalves and Sanchez
- Space as an organising dimension: Belgum, Handley and Bott
- Workflows

Space and digital scholarship

- Space and place can be a greater part of scholarship because:
 - Much easier to incorporate images, maps etc.
 - GIS allows geographic analysis of a wide range of data



Geographic Information Systems (GIS)

- Systems which allow analysis of geocoded data
- For example:
 - Census data gives demographic details
 - Add geocoding at desired level of detail
 - Analyses can be carried out by geographic location
- Very easy to present results as a map or maps
 - But this is not essential for GIS analysis



GIS in the three articles

- Thomas and Ayers took large amounts of historical data and made it usable for GIS analysis
- Goncalves and Sanchez used data which came already geocoded
 - Carried out GIS-style analysis on it
- Belgum et al use geoparsing to make data available for geographic analysis
 - Their end product is a hand-tooled GIS system

Thomas and Ayers – what they did

- Look at two communities which are as similar as possible
- Except for presence/absence of slavery
- Chosen communities had to be:
 - Close together (both in Shenandoah Valley region)
 - Similar economic possibilities



Thomas and Ayers: data sources

- Article is based on a rich array of data:
 - Text sources newspapers, diaries, letters etc.
 - Census records
 - Contemporary maps
 - (others)

• Data made digital and transformed as necessary for GIS analysis

		VINDICAT LADVERTISER.	
	AND GENERAL	THE CONSTITUTION.	To be discharged by Two Dollars to Advance
TERMS, \$2.50 Cents,	THE UNION, BASED UP	INTY. VIRGINIA, JUNE 16, 1859.	NUMBER 24,
VOLUME XV.	STAUNTON, AUGUSTA COC AL. MISCELLANEOUS. Sda	the Beilin	ATTIMES FROM OUR DIFLOMATE IN a -The following is telegraphed from ington to the New York Herald: BALTIMORE, April 50th; 1859.
The Bindicator, DR. JOHNSTON	GREAT ATTRACTIOS & BARGAINS	TLE, LOVE ME LONG." a voice behind us; "you will remain a good-"D. for-nothing fellow all your life." State	repartment, brought by the Asia. Dear Sir: - We have frequently and the Department, brought by the Asia our Ministers at Loudon, Brance and our circulars, and have often have the

dh@monash - ATS1208

Thomas and Ayers: access to data

- Selected data directly supports argument
 - Available through Points of Analysis tab
 - More data is available on associated website
- Underlying data is much more accessible than in conventional (printed) work
- A lot of the data is about place and space



Thomas and Ayers: presentation

- The authors set out to be innovative
- Features:
 - Making data easy to reach
 - Allowing reader to choose a path through material
 - No print version exists
 - The physical journal has an overview only (8 pages long)
 - No downloadable pdf or anything similar
- Radical approach, still impressive after 15 years

Goncalves and Sanchez: what they did

- Data is tweets in Spanish
- All geolocated tweets in Spanish collected for 2 years
- Total dataset $> 5 \times 10^7$ tweets
- Built a list of concepts each represented by several possible words
- 7.5 x 10⁵ tweets in the full dataset contained one of these words



Goncalves and Sanchez: the maths

- Geographic division into cells 0.25° x 0.25° (c 25km² at the equator)
- Dominant word for each concept in each cell calculated by simple majority rule
- Results in matrix of 1135 rows (geographic cells) and 131 columns (concepts)
- Each point in matrix has value of 1 or 0
- Example of what it looks like on next slide

Cars and computers



Machine learning

- Clustering of data was done with machine learning method
- Two steps:
 - Principal Components Analysis to reduce dimensionality (first 40 components included 94% of the variance)
 - Clustering with K-means algorithm
 - Best solution is 2 clusters

What it all means

- First look at geographic make up of clusters suggested population density was relevant
- Data on population density for each geographic cell imported from an existing dataset
 - LandScan 2007 High Resolution global Population Data Set
 - Using a GIS data source
- Next slide shows boxplot comparing population density for two clusters

Superdialects and population



What it all means

- Standard story has been that major split in Spanish is between European variety and New World variety
- Goncalves and Sanchez show that major split is between an international urban variety and rural varieties
 - Cluster β can be split into several regional varieties
- Caveats:
 - Based on lexical variation pronunciation variation is also important
 - Population may not be representative (younger, good knowledge of technology,....)

The Differences Between

European and Latin American

Spanish

thewanderinglinguist.com

Belgum, Handley and Bott: what they did



- Starting point is bibliographies of travel books published in C19
 - 3000 works published in Britain
 - Titles taken as basic data
- Each work was located geographically on the basis of places mentioned in the title
- Results can be accessed via interactive map

Distant reading

- Idea introduced by Franco Moretti
- Opposed to traditional literary analysis: close reading
- Moretti published an essay on the titles of English novels in C18 and C19
- He correlates changes and linguistic patterns with changes in society and publishing
- Current study is less ambitious first question is which places were written about (and when)



First attempt

- Used bibliography of French travel writing
- Place names identified by hand (student assistants)
- Relied on coder knowledge and reference works
- Processing 100 entries took on average one hour
 - Plus time to review and correct errors

Geoparsing

- Takes unstructured references to locations
- Tries to resolve them to unambiguous identifiers (e.g. co-ordinates)
- Steps:
 - Identifying references in text (cf. Named Entity Recognition)
 - Match text reference to entry or entries in geographic database
 - Resolve ambiguities if possible
- Problems:
 - Variant names (archaic names and spellings)
 - Ambiguities Paris, France v. Paris, Texas

Geographic database



- GeoNames is a very large digital gazetteer
- Stores variant names
- But no information on frequency of use
- This project used Wikipedia
 - Entries on places typically include geographic identifier (latitude and longitude)
 - Good coverage of higher level topography
 - Not many travel books have e.g. village name in title
 - Search history can be used for disambiguation
 - Paris, France is a common search term, Paris, Texas is not

Searching Wikipedia

- Procedure started with 5 word strings from title, then 4 word strings and so on down to individual words
 - Aim was to match e.g. New York City before New York
- Match of string to Wikipedia entry with geographic coordinates taken as identification of a place name
- Place name and co-ordinates added to bibliography database
- Manual checking as final step, with a map interface

Image: constraint of the second se	Location	of Catania	[show]		
Country Italy Region Sicily Metropolitan Catania (CT) city Bicocca, Codavolpe, Junghetto, Pantano d'Arci, Paradiso degli Aranci, Passo Cavaliere, Passo del Fico, Passo Martino, Primosole, Reitano, Vaccarizzo, Villaggio Delfino Government Nacorizzo, Villaggio Delfino • Mayor Salvo Pogliese (FI) Area[1] . • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01)[2] . • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 tg² Patron saint St. Agatha Saint day February 5 Nucherin C	Location of Catania in Sicily				
CountryItalyRegionSicilyMetropolitanCatania (CT)cityBicocca, Codavolpe,Junghetto, Pantano d'Arci, Paradiso degli Aranci, Passo Cavaliere, Passo del Fico, Passo Martino, Primosole, Reitano, Vaccarizzo, Villaggio DelfinoGovernmentNayorSalvo Pogliese (FI)Area[1] • Total182.9 km² (70.6 sq mi)Elevation7 m (23 ft)Population (2018-01-01)[2] • Total311,620• Density1,700/km² (4,400/sq mi)Demonym(s)CataneseTime zone • Summer (DST)UTC+2 (CEST)Postal code Dialing code995ISTAT code087015 tg²Patron saint Saint daySt. Agatha Saint day	Coordinates: 🜉 37°30'0"N 15°5'25"E				
Netropolitan city Stary Frazioni Bicocca, Codavolpe, Junghetto, Pantano d'Arci, Paradiso degli Aranci, Passo Cavaliere, Passo darlino, Primosole, Reitano, Vaccarizzo, Villaggio Delfino Government Navor • Mayor Salvo Pogliese (FI) Area[¹¹] . • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] . • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015tg² Patron saint St. Agatha Saint day February 5 Start code	Country	Italy			
city city Frazioni Frazioni Bicocca, Codavolpe, Junghetto, Pantano d'Arci, Paradiso degli Aranci, Passo Cavaliere, Passo darino, Primosole, Reitano, Vaccarizzo, Villaggio Delfino Government • Mayor Salvo Pogliese (FI) Area[¹¹ • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015tg² Patron saint St. Agatha Saint day February 5	Metropolitan	Catania (CT)			
Frazioni Bicocca, Codavolpe, Junghetto, Pantano d'Arci, Paradiso degli Aranci, Passo Cavaliere, Passo del Fico, Passo Martino, Primosole, Reitano, Vaccarizzo, Villaggio Delfino Government • Mayor Salvo Pogliese (FI) Area[¹¹ 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95 ISTAT code 087015tg² Patron saint St. Agatha Saint day Varbaits Official unschain t ^C	city				
Government • Mayor Salvo Pogliese (FI) Area[1] . • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01)[2] . • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 ts² Patron saint St. Agatha Saint day February 5	Frazioni	Bicocca, Codavolpe Junghetto, Pantano Paradiso degli Aran Cavaliere, Passo de Passo Martino, Prim Reitano, Vaccarizzo Villagoio Delfino	e, o d'Arci, ici, Passo el Fico, nosole, o,		
• Mayor Salvo Pogliese (FI) Area ^[1] • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015tg² Patron saint St. Agatha Saint day February 5	Government				
Area[1] • Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] • • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 08701519 Patron saint St. Agatha Saint day February 5	Mayor	Salvo Pogliese (FI)			
• Total 182.9 km² (70.6 sq mi) Elevation 7 m (23 ft) Population (2018-01-01) ^[2] . • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 ts² Patron saint St. Agatha Saint day February 5	Area ^[1]				
Elevation 7 m (23 ft) Population (2018-01-01) ^[2] • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 t2 ⁹ Patron saint St. Agatha Saint day February 5	Total	182.9 km ² (70.6 sq	mi)		
Population (2018-01-01) ^[2] • Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 087015 tg? Patron saint St. Agatha Saint day February 5	Elevation	7 m (23 ft)			
• Total 311,620 • Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 t2 Patron saint St. Agatha Saint day February 5	Population (2018-0)1-01) ^[2]			
• Density 1,700/km² (4,400/sq mi) Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 t2 Patron saint St. Agatha Saint day February 5	Total	311,620			
Demonym(s) Catanese Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 @ Patron saint St. Agatha Saint day February 5	Density	1,700/km ² (4,400/s	q mi)		
Time zone UTC+1 (CET) • Summer (DST) UTC+2 (CEST) Postal code 95100 Dialing code 095 ISTAT code 087015 @ Patron saint St. Agatha Saint day February 5 Website Offsic/ussbatier C	Demonym(s)	Catanese			
Postal code 95100 Dialing code 095 ISTAT code 087015 @ Patron saint St. Agatha Saint day February 5	Time zone • Summer (DST)	UTC+1 (CET) UTC+2 (CEST)			
Dialing code 095 ISTAT code 087015 @ Patron saint St. Agatha Saint day February 5	Postal code	95100			
ISTAT code 087015129 Patron saint St. Agatha Saint day February 5	Dialing code	095			
Patron saint St. Agatha Saint day February 5	ISTAT code	087015 🗗			
Website Official website C	Patron saint Saint day	St. Agatha February 5			
Website Official Website	Website	Official website @			

Remaining problems

- Small number of cases where Wikipedia information did not align well with book
 - Manual input required
- Assigning point locations to regions is a general problem
 - Broad regions in titles (e.g. *The Islands of Greece*) were hard to identify

Results

- Online tool to visualise the data
- Provides details of individual works
- Possibility to view variation over time
- Extensible system more data can be added

Workflows

- Each article details methods
- In each case, initial data needed to be manipulated in some way for further use
- Key element in digital humanities scholarship:
 - Finding structure in fully/partially unstructured data
 - Transforming data to make structure clear and available for analysis

What we spend our time on.....



Belinda Weaver

Data preparation, including clean up, is 80% of the data work in research, according to a survey forbes.com/sites /gilpress...



References

- Belgum, Kirsten, Keith Handley & Rachel Bott. 2018. Mapping travel writing: a digital humanities project to visualise change in nineteenth-century published travel texts. *Studies in Travel Writing* 22(3). 306–324. doi:10.1080/13645145.2019.1575765.
- Gonçalves, Bruno & David Sánchez. 2014. Crowdsourcing Dialect Characterization through Twitter. (Ed.) Tobias Preis. *PLoS ONE* 9(11). e112074. doi:10.1371/journal.pone.0112074.
- Thomas III, William G. & Edward L. Ayers. 2003. The differences slavery made: A close analysis of two American communities. *The American Historical Review*. <u>http://www2.vcdh.virginia.edu/AHR/</u>.