



Beatrice Alex balex@inf.ed.ac.uk

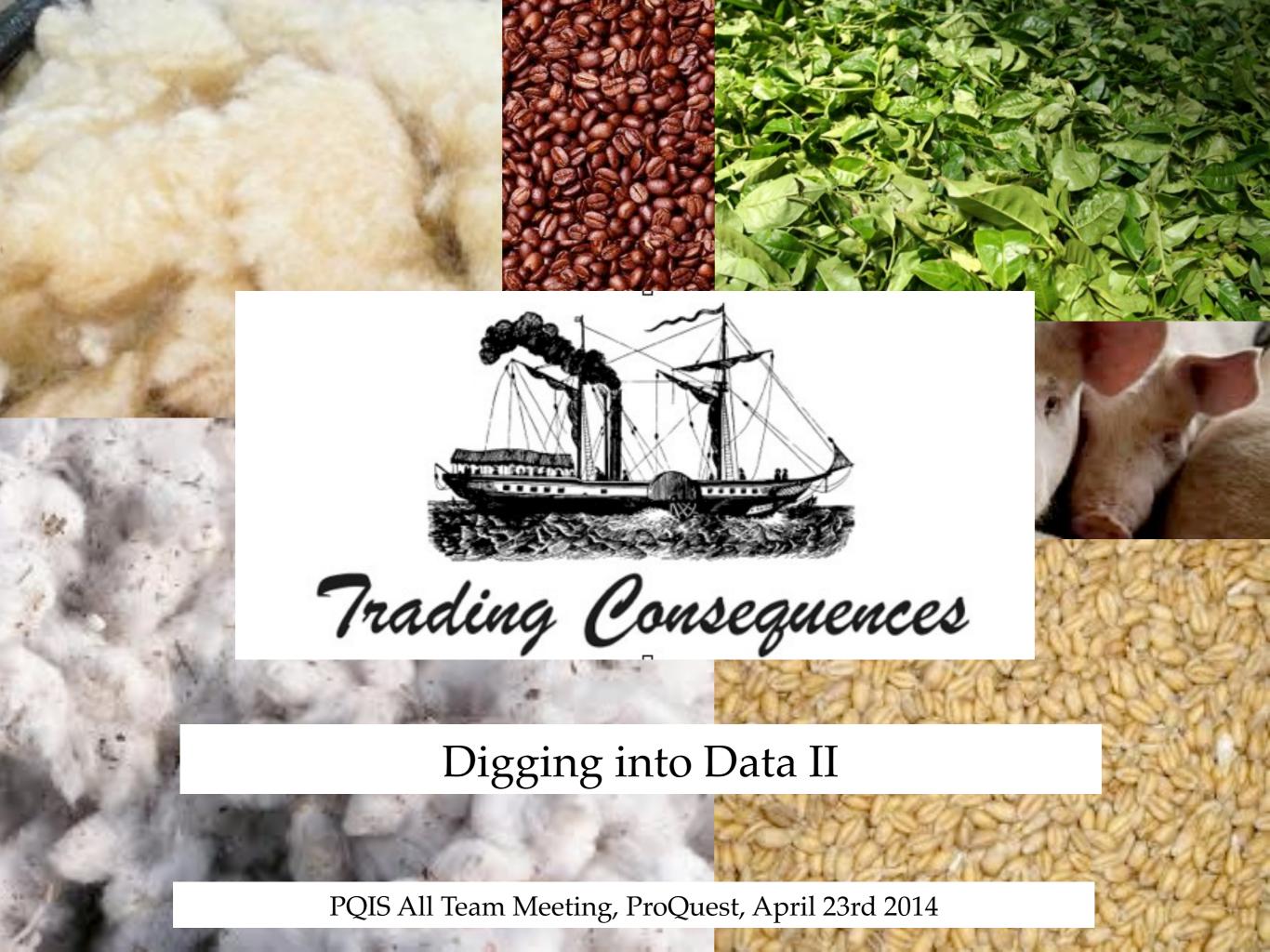


TEXT MINING

- Describes a set of linguistic, statistical and/or machine learning techniques that model and structure the information content of textual resources.
- Turns unstructured text into structured data (e.g. relational database or linked data).
- Is very useful for analysing large text collections automatically (overcoming data paralysis).
- Goal in DHSS research: By analysing large amounts of textual data, help HSS scholars to discover novel patterns and explore hypotheses.

TYPES OF ANALYSES

- Named entity recognition.
- Grounding, e.g. geo-referencing.
- Relation extraction.
- Clustering, e.g. topic modelling.
- Sentiment analysis.



PROJECT TEAM







Colin Coates, Andrew Watson: historical analysis



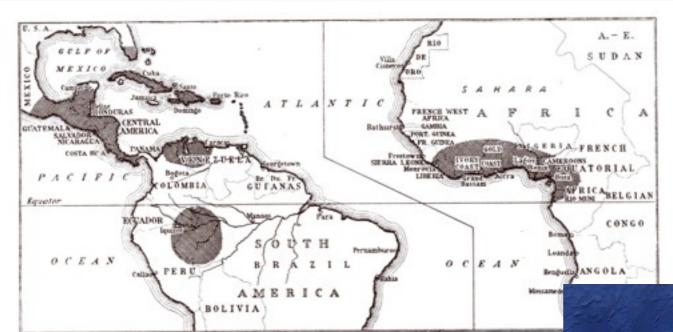
Jim Clifford: historical analysis



James Reid, Nicola Osborne: data management, social media

Aaron Quigley, Uta Hinrichs: information visualisation

TRADITIONAL HISTORICAL RESEARCH



Map showing the areas where mahogany is grown

Gillow and the Use of Mahogany in the Eighteenth Century, Adam Bowett, Regional Furniture, v.XII, 1998.



PQIS All Team Meeting, ProQuest, April 23rd 2014

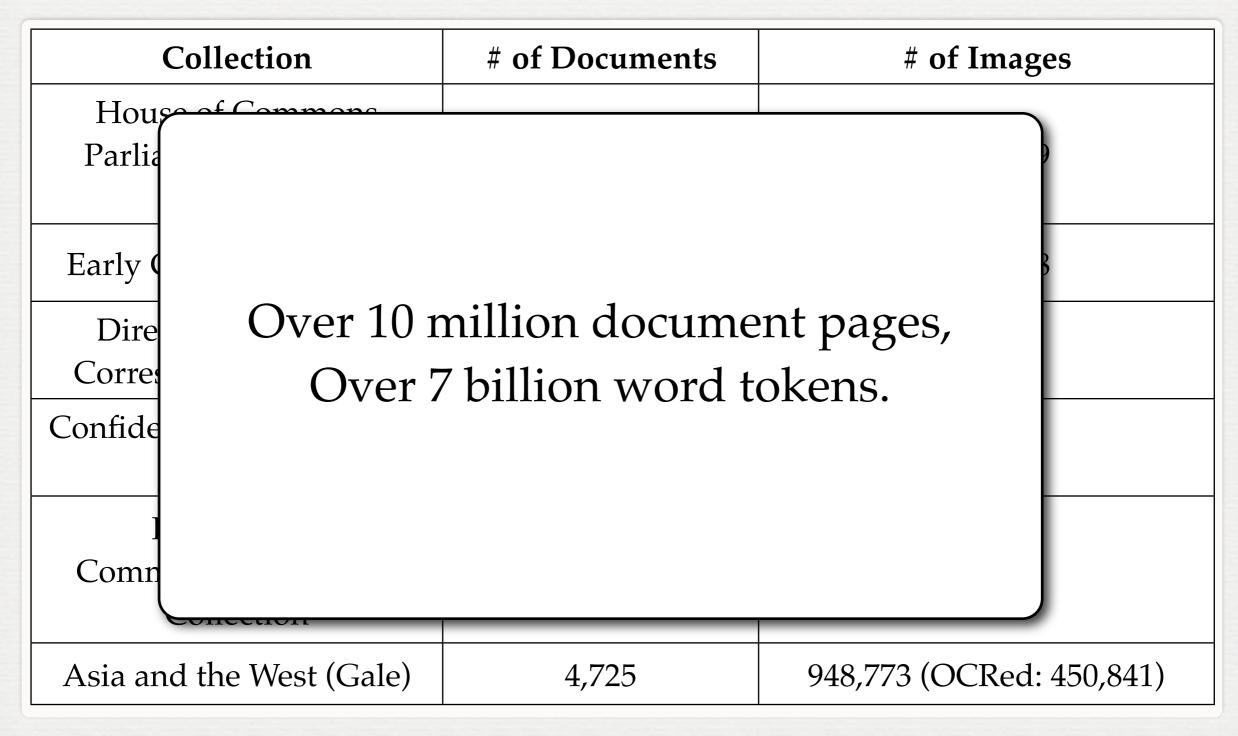
PROJECT GOALS

- Text mining, data extraction and information visualisation to explore big historical datasets.
- Focus on how commodities were traded across the globe in the 19th century.
- Help historians to discover novel patterns and explore new research questions.

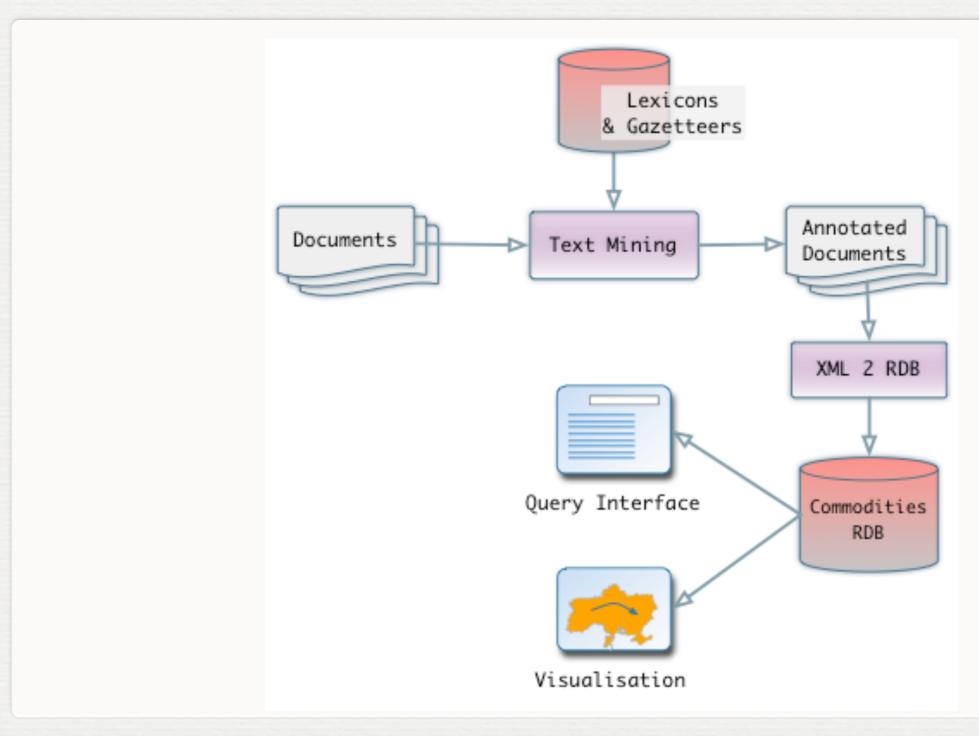
DOCUMENT COLLECTIONS

Collection	# of Documents	# of Images
House of Commons Parliamentary Papers (ProQuest)	118,526	6,448,739
Early Canadiana Online	83,016	3,938,758
Directors' Letters of Correspondence (Kew)	14,340	n/a
Confidential Prints (Adam Matthews)	1,315	140,010
Foreign and Commonwealth Office Collection	1,000	41,611
Asia and the West (Gale)	4,725	948,773 (OCRed: 450,841)

DOCUMENT COLLECTIONS



ARCHITECTURE



MINED INFORMATION

Example sentence:

From Padang was exported, in 1871, 6,127 piculs of cassia bark, of which a large portion was shipped to America (Fliickiger and Hanbury). ...

- Normalised and grounded entities:
 - commodity: cassia bark [concept: Cinnamomum cassia]
 - date: 1871 (year=1871)
 - location: Padang (lat=-0.94924;long=100.35427;country=ID)
 - location: America (lat=39.76;long=-98.50;country=n/a)
 - quantity + unit: 6,127 piculs

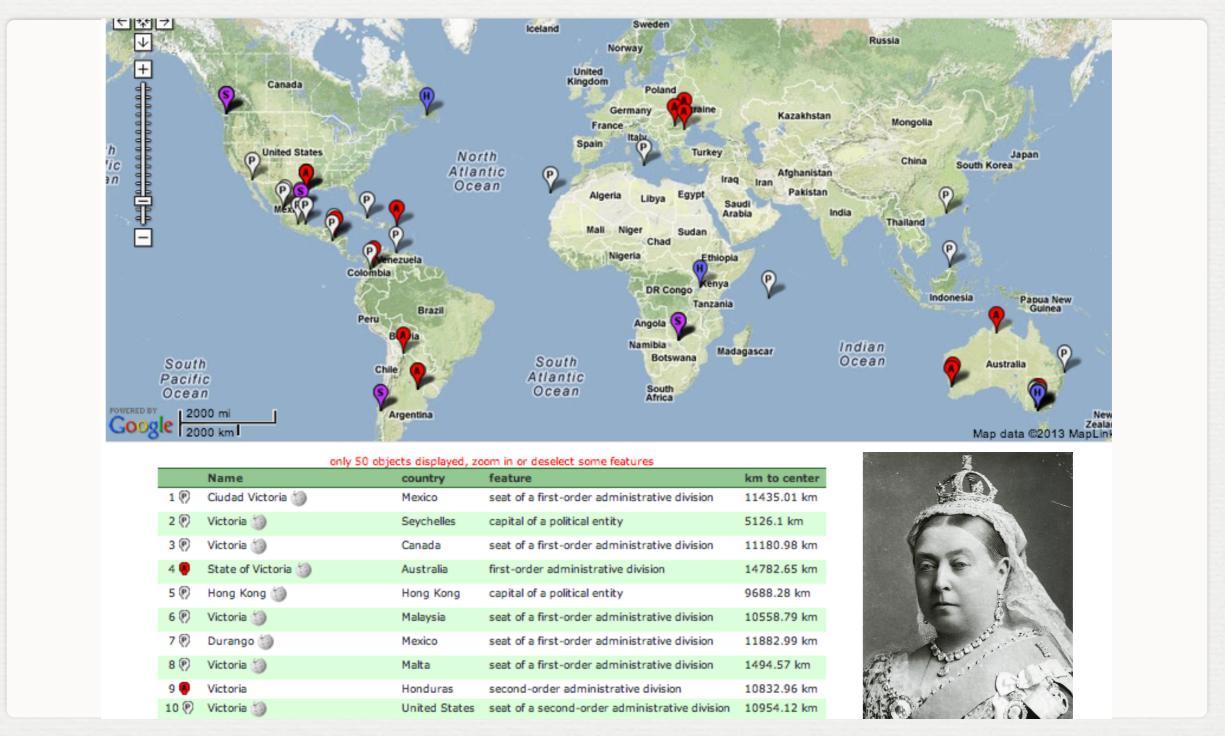
MINED INFORMATION

Example sentence:

From Padang was exported, in 1871, 6,127 piculs of cassia bark, of which a large portion was shipped to America (Fliickiger and Hanbury). ...

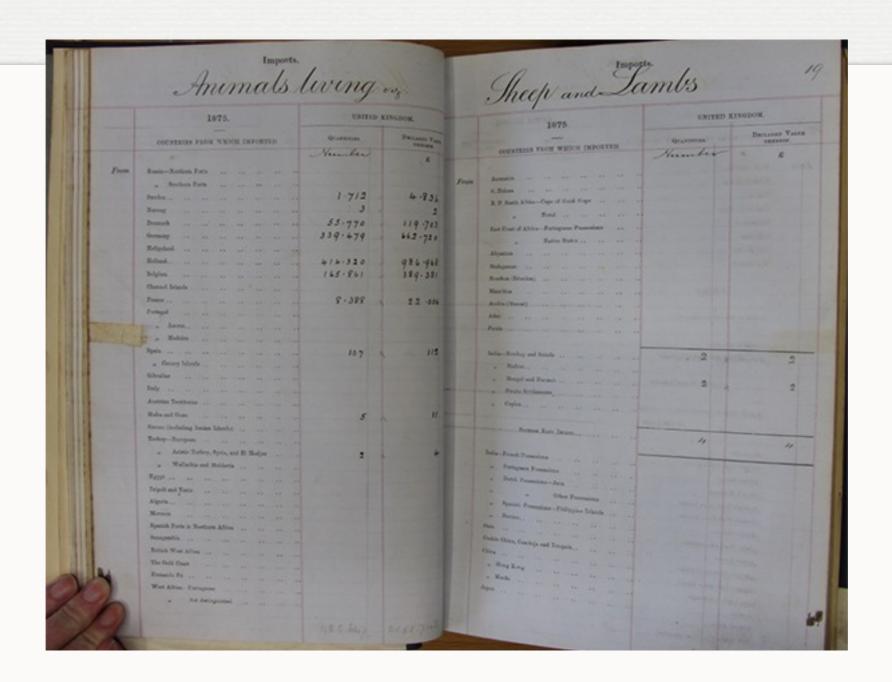
- Extracted entity attributes and relations:
 - origin location: Padang
 - destination location: America
 - commodity–date relation: cassia bark 1871
 - commodity-location relation: cassia bark Padang
 - commodity-location relation: cassia bark America

EDINBURGH GEOPARSER



PQIS All Team Meeting, ProQuest, April 23rd 2014

COMMODITY LEXICON



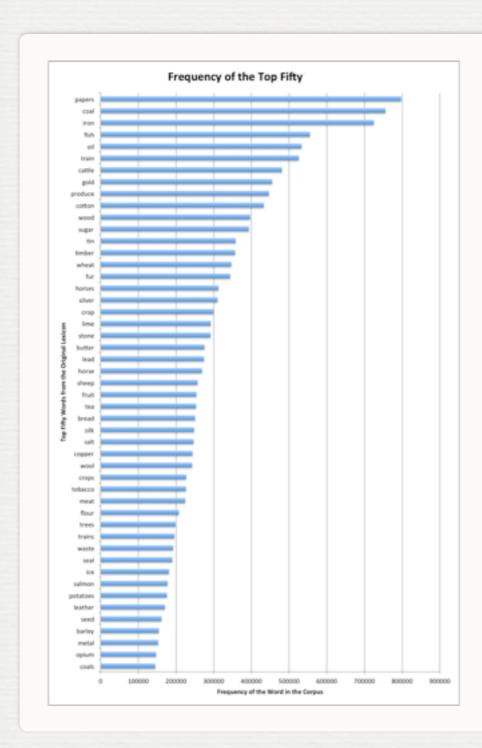
Seed set from customs import records.

LEXICON CREATION

Seed lexicon	~600
Extended lexicon	~17,000
With pluralisation of single word entries	~20,500

Bootstrapping a historical commodities lexicon with SKOS and DBpedia. Klein, Alex & Clifford, LaTeCH 2014.

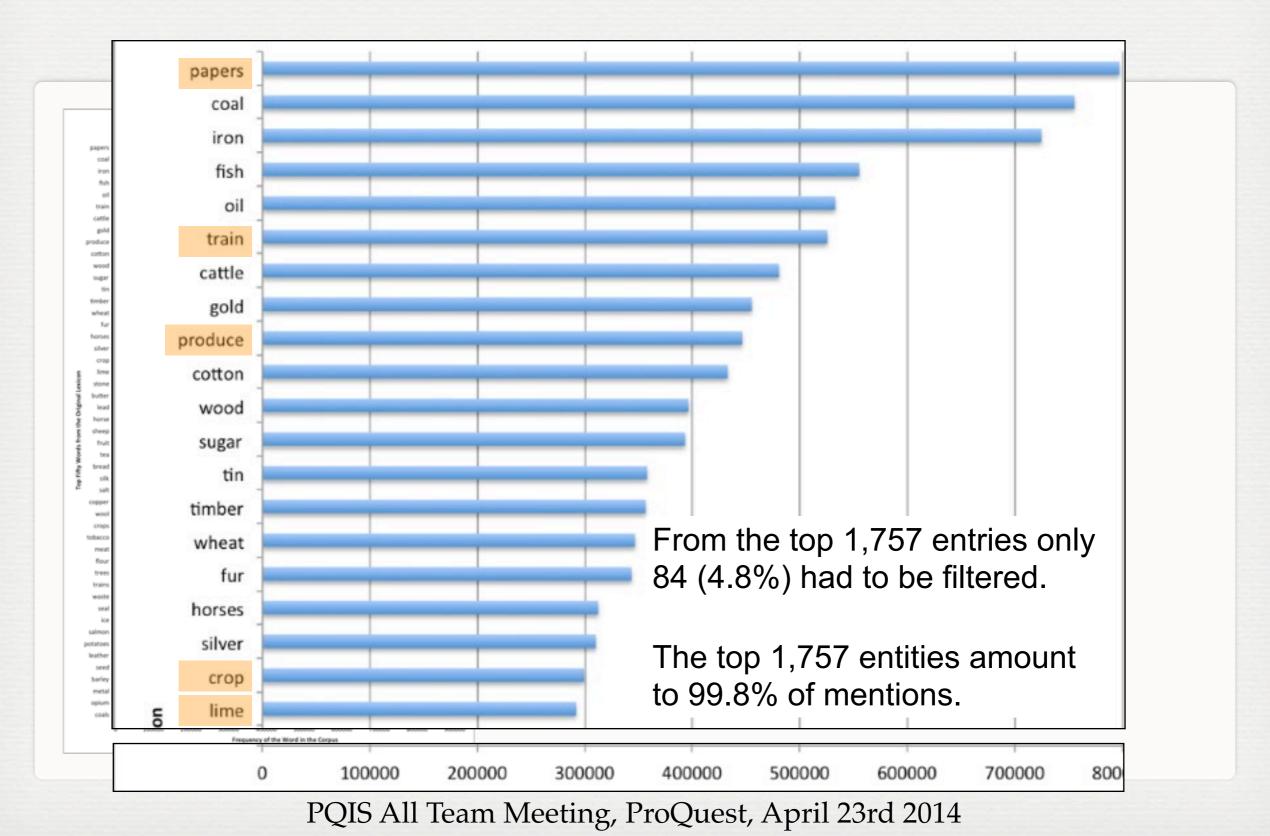
LEXICON CLEAN-UP



From the top 1,757 entries only 84 (4.8%) had to be filtered.

The top 1,757 entities amount to 99.8% of mentions.

LEXICON CLEAN-UP



NOISY DATA

- Optical character recognition contains many errors and often the structure of the page layout is lost.
 - Sophistication of the OCR engine and scanning equipment.
 - Quality of the original print and paper.
 - Use of historical language.
 - Information in page margins (header, page numbers, etc.).
 - Information in tables.
 - Language of the text.

FIXING NOISY DATA

- Text normalisation and correction:
- End-of-line soft hyphen removal
 - Dehyphen all token-splitting hyphens using a dictionary-based approach.
- "False f"-to-s conversion
 - Convert all false f characters to s using a corpus.
- Example: reduced number of words unrecognised by spell checker from 61 to 21 -> 67%, on average 12% reduction in word error rate in a random sample (Alex et al., 2012).

FIXING NOISY DATA

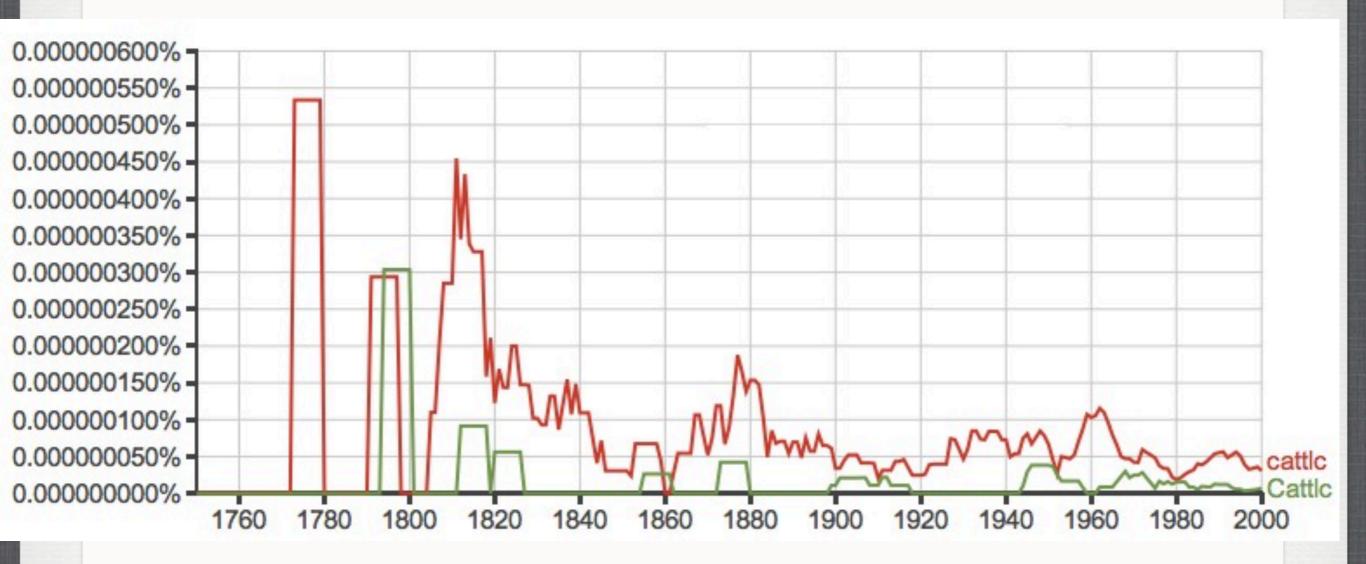
```
<?xml version="1.0" encoding="UTF-8"?>
<article id="10.2307/60227644">
</pr
MR. WILLIAM BELSHAM. BY HERBERT MARSH, B. D. F. R. S. and tellow or st. John's college, Cambridge. Ecntmn:
PRINTED FOR JOHN STOCKDALE, PICCADILLY. 1801. t35)lvjf~ Udf4~ P.]]> </page>
<page> <![CDATA[T, G1IXET, Printer.]]> </page>
(CDATA[' INTRODUCTION, AS the following Vindication may fall into the hands of perfons who have
never read the Hiftory of the Politicks of Great Britain and France, it will not be improper, before I enter
on my Defence, to ftate the principal facts, which were fuccef- fively proved by authentic documents, in the
fixteen chapters, of which that wrork is compofed. - r 1. In the celebrated conference at PiUnitz; in Auguft,
i;gi, the Britifh Government took not the rnoft diftant part: and if.-any treaty was concluded there, which
is itfelf a matter of great doubt, the Britifh Go- vernment not only never acceded to it, but was, never
apprifed even of its contents.- Further, when the Britiih Government was requefted in 1701 to join a
coalition againft France, it gave a pofitive and unequivocal refufal. B 2 2. Toward]]> </page>
<page> <![CDATA[4 2. Toward the clofe of the fame year the valuable colony of St. Domingo was pre- served to
France by the timely affiftance fent by Lord Effingham, then Governor of Jamaica : and the British. Cabinet
fignified through its AmbafTador at Paris to the French Government, that it fully approved of Lord
Effingham's conduct.. At the fame time, true to the ftri&eir. principles of ho- nour and neutrality, it
refufed the advan- tageous offer made by the French colonifts, who were highly diflatisfied with the Na-
tional AfTembly, to furrender the French part of St. Domingo to the Crown of Bri- tain. And thefe a6ls of
generofity were re* paid by France with the utmoft ingrati- tude. 3. When Louis XVI. formally accepted the
new conflitution, in September, 17Q1, and fent circular letters to the different Courts of Europe fignifying
his affent, the Court of Great Britain was one of the firft which returned an anfwer; and the anfwer was
couched in very refpectful terms, where- as fome other courts either did not anfwer at HfWta]]> </page>
</article>
                         PQIS All Team Meeting, ProQuest, April 23rd 2014
```

FIXING NOISY DATA

```
<document>
<meta>
<attr name="docid">10.2307/60227644</attr>
</meta>
<text>
THE HISTORY OP THE POLITICKS OF GREAT BRITAIN AND FRANCE, VINDI GATED FROM A LATE ATTACK OF MR. WILLIAM
BELSHAM. BY HERBERT MARSH, B. D. F. R. S. and tellow or st. John's college, Cambridge. Ecntmn: PRINTED FOR
JOHN STOCKDALE, PICCADILLY. 1801. t35)lvjf~ Udf4~ P.
T, G1IXET, Printer.
' INTRODUCTION, AS the following Vindication may fall into the hands of persons who have never read the
History of the Politicks of Great Britain and France, it will not be improper, before I enter on my Defence,
to state the principal facts, which were successively proved by authentic documents, in the sixteen chapters,
of which that wrork is composed. r 1. In the celebrated conference at PiUnitz; in August, i;gi, the British
Government took not the rnoft distant part: and if.-any treaty was concluded there, which is itself a matter
of great doubt, the British Government not only never acceded to it, but was, never apprised even of its
contents.- Further, when the Britiih Government was requested in 1701 to join a coalition against France, it
gave a positive and unequivocal refusal. B 2 2. Toward
4 2. Toward the close of the same year the valuable colony of St. Domingo was preserved to France by the
timely assistance sent by Lord Effingham, then Governor of Jamaica : and the British. Cabinet signified
through its AmbafTador at Paris to the French Government, that it fully approved of Lord Effingham's
conduct.. At the same time, true to the ftri& eir. principles of honour and neutrality, it refused the
advantageous offer made by the French colonists, who were highly diflatisfied with the National AfTembly, to
surrender the French part of St. Domingo to the Crown of Britain. And these a6ls of generosity were re* paid
by France with the utmost ingratitude. 3. When Louis XVI. formally accepted the new conflitution, in
September, 17Q1, and sent circular letters to the different Courts of Europe signifying his assent, the Court
of Great Britain was one of the first which returned an answer; and the answer was couched in very
refpectful terms, whereas some other courts either did not answer at HfWta
</text>
                           PQIS All Team Meeting, ProQuest, April 23rd 2014
</document>
```

```
Proclamations, Pro * v mie RL' E.LI S B.AIGO7.
iVICTORIFIA. h> I 1/ t(aT'' of' GO!>. tif ih Firi.
ea fil~T/ r<' lluil'tIT, (i'. i', QUEE'. Tc, iii-
n iTiV i ' ui tillhT'nt, 111te 1 eihT' Colin. ('it; ZI-s.
uni 14Lt1ussuls ce t rib ev iii tJ1u stat. it have Iei
t's.iiititTud ztntt liild, a tutt &lt; A 11i10C.
```

Extract of Early Canadiana Online document 9_00952_3, p. vi.





Google books Ngram Viewer



HOW NOISY IS TOO NOISY?

qBiu si }S3A:req s,uauuaqsu aq} }Bq} uirepo.ifT 'papua}X3 sSuiav }qSuq Jiaq} qiiM jib ui snnS bbs aqx 'a"3(s aq} tnojj ssfitns q}TM Sni5[ooi si jb}s }S.ii; aqx 'papnaoSB q}Bq naABSjj qS;H °1 ssbui s.uauuaqsu aqx

Extract from document 10.2307/60238580 in FCOC.

HOW NOISY IS TOO NOISY?

The fishermen's mass to High Heaven bath ascended,

The first star is looking with smiles from the sky,

The sea gulls in air with their bright wings extended,

Proclaim that the fishermen's harvest is nigh.

Oh! I have been there, and seen them all kneeling,
Their faces turned upward beneath the bright sky,
And heard their low throbs of worshipping feeling,
And seen pure devotion gash out from each eye;
And oft, indd the worlds rough battles and sorrow,
I've thought of that fishermen's mass of the sen,
And pray'd for such faith in brighter to-morrow
As shone from each face in the bay of Tralee.

By fair Trales bay, oh! did you e'er wander,
Where quiet and beauty are blended in one,
Watching the brooklets their pearly drops squander,
Bright gens soon to form for the brow of the sun?
And have you e'er heard the swell of devotion,
Like chorus of angels from depths of the sea,
As mass hath been said by those sons of the ocean,
The fishermen humble of quiet Trales?

The Fishermen's Mass.

Stronger than the tempest's rage is Thy power from age to age;

Is Thy power from age to age;

Thou didst once its wrath assuage,

We will own Thy mighty power,

Even in the darkest hour,

And receive as blessed dower

And Thy will,

When upon the stormy deep,

Fishermen their vigils keep,

Wherefore should their lovers weep?

When the storm is raging high,
And the waves leap to the sky,

Wherefore should we fret and sigh?

Thou dost eare.

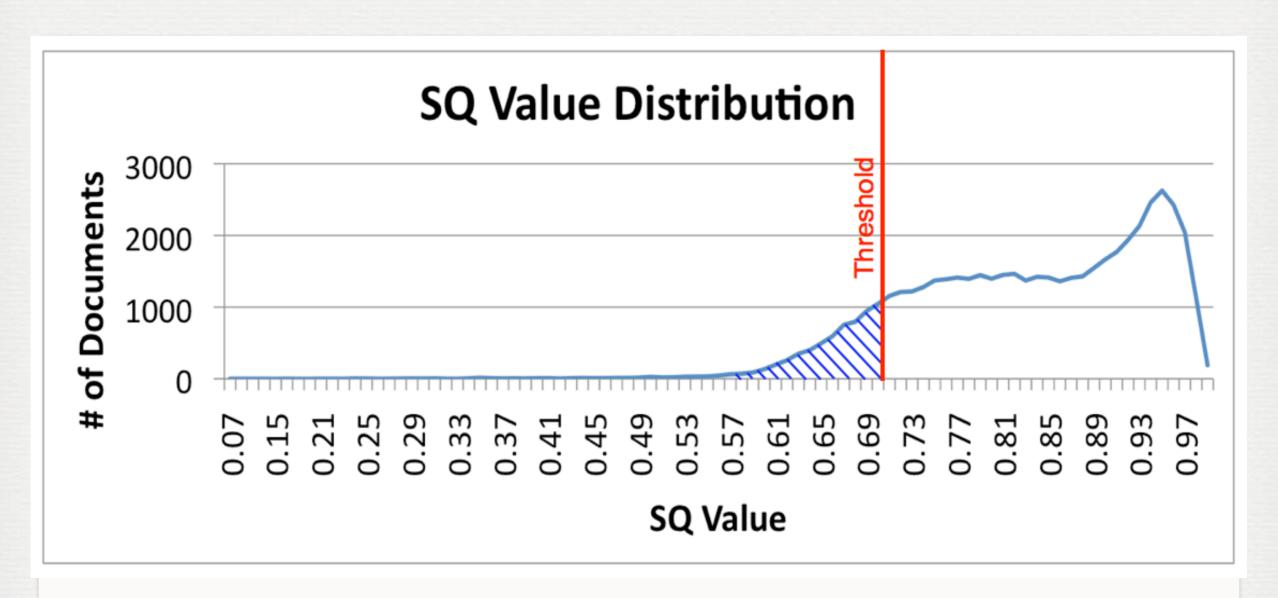
Again and again with fond kisses bedewing

Her baby's soft check, as the rosebud as fair;
Again and again her look upward ronewing,
This song, sweet and solemn, she breathes on the air:

She looks through her tenrs to her Father in heaven, Then, kissing her child, cries, "Ah, bless'd be that smile! And bless'd be the saints! now the promise is given! Sure, mass it was said—he'll he here in a while."

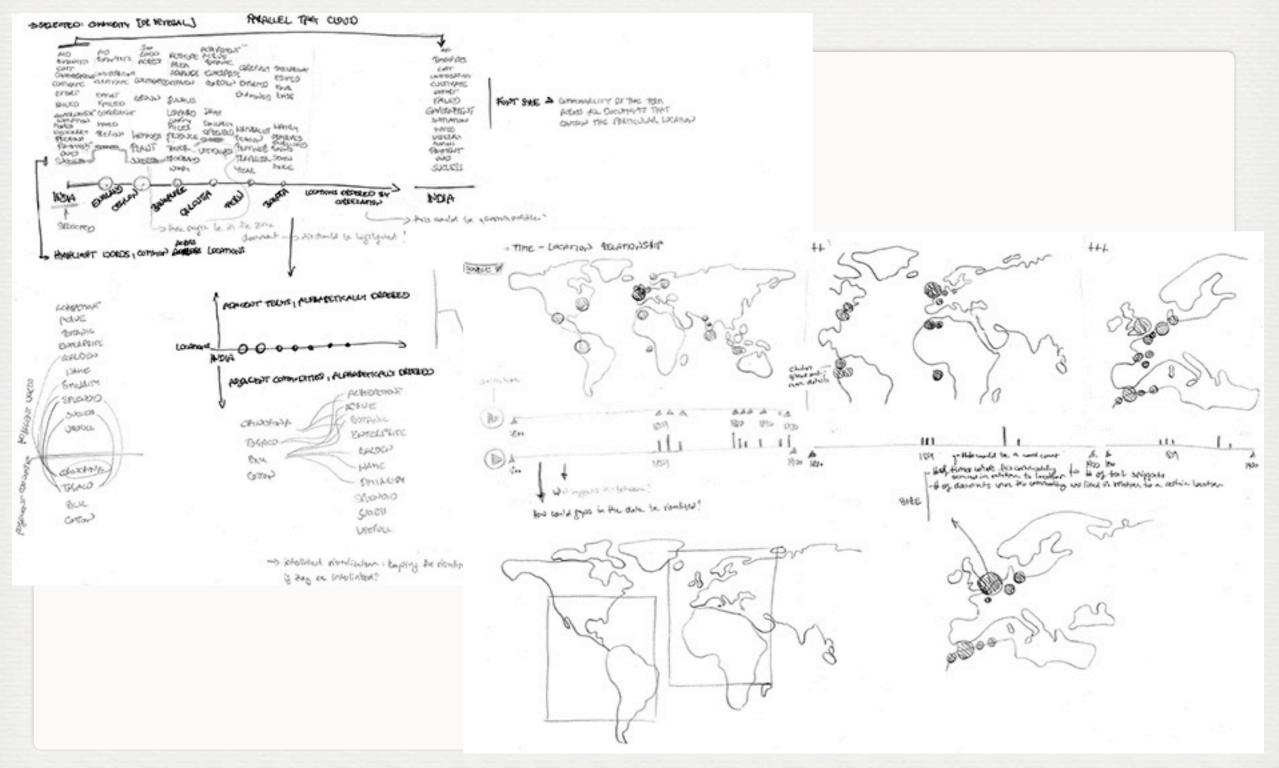
Upon the dark storm, from the fisherman's dwelling, There shines a faint light like a star through a cloud; Fast bugging her babe to her bosom high swelling, Tight clasping her rosary, Mary's low bowed.

Brave Dermot looks up with a spirit all lowly, And, wiping the brine from his storm-besten head, Cries " Mother of Jesus, and saints the most holy, Betriend!—Sure, the fishermen's mass it was said."

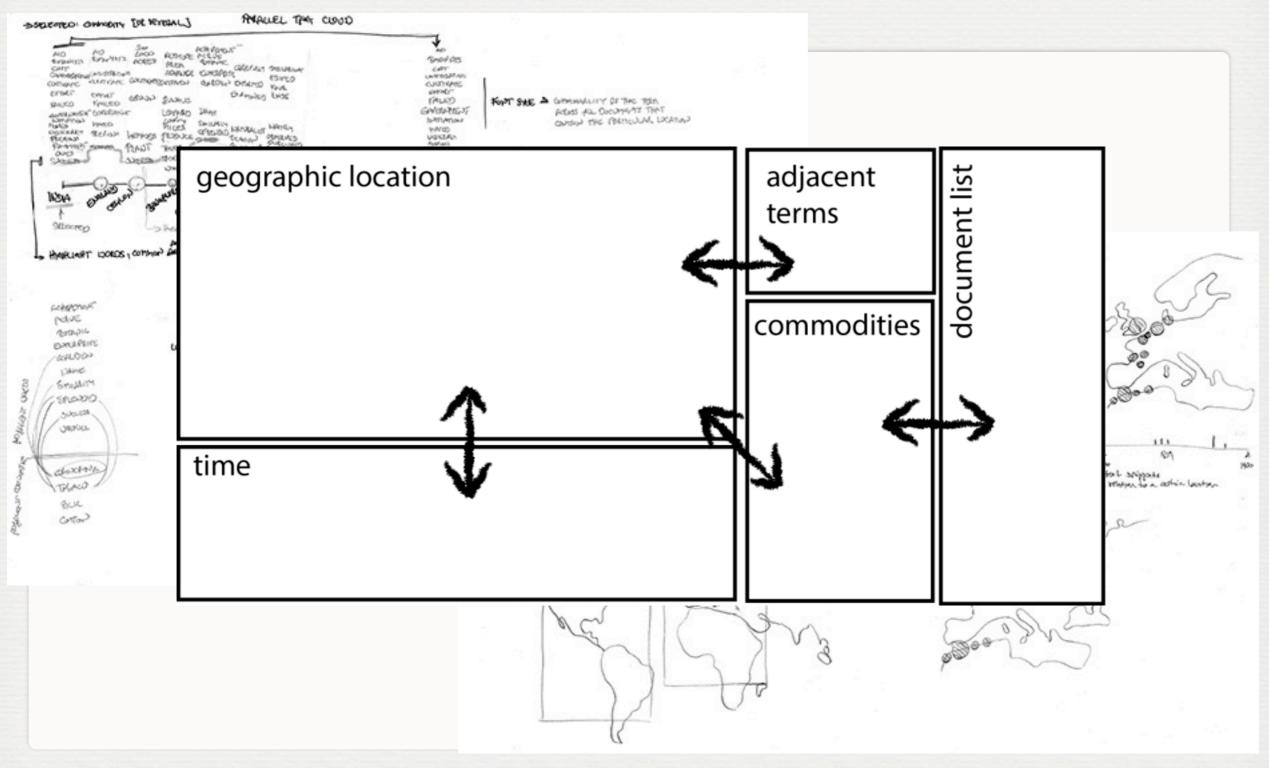


Study of correlating manual quality ratings of documents with automatic quality scoring (Alex & Burns, DATeCH 2014).

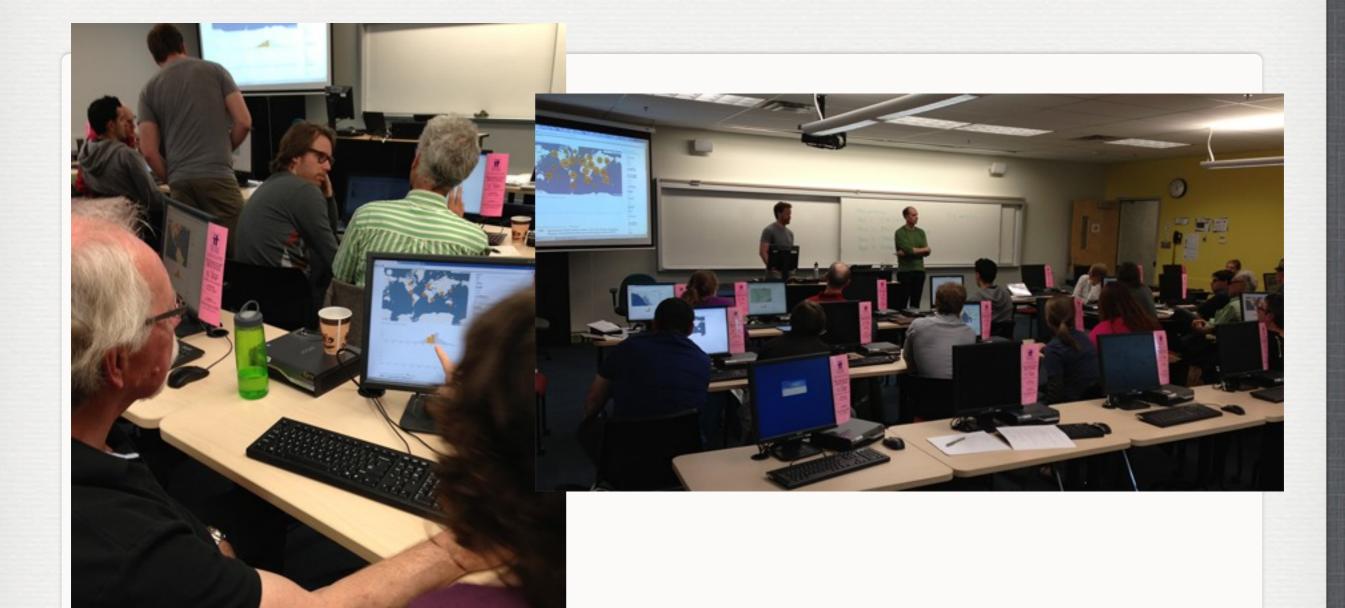
VISUALISATION SKETCHES



VISUALISATION SKETCHES

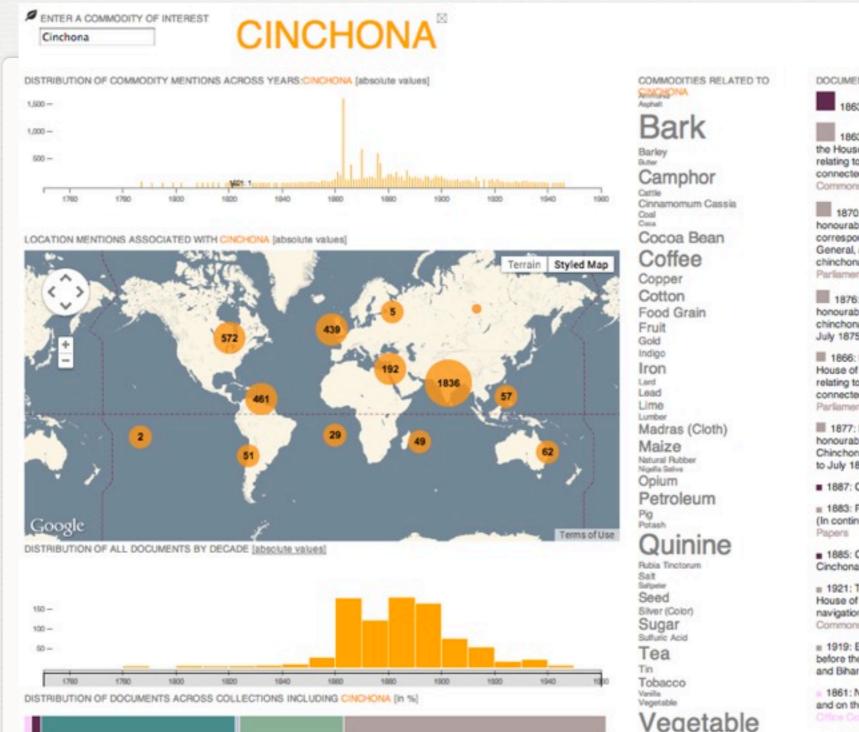


USER WORKSHOP



User workshop to improve the functionality of the interface (Hinrichs et al., 2014)

BRINGING ARCHIVES ALIVE



Vegetable Ivory



DOCUMENTS INCLUDING CINCHONA [top 100]

1863: East India (Chinchona Plant); Miscellaneous

1863: East India (chinchona plant). Return to an address of the Honourable the House of Commons, dated 9 March 1863;--for, "copy of correspondence relating to the introduction of the chinchona plant into India, and to proceedings connected with its cultivation, from March 1852 to March 1863."; House of Commons Parliamentary Papers

1870: East India (chinchona cultivation). Return to an address of the honourable the House of Commons, dated 3 May 1870;—for "copy of all correspondence between the Secretary of State for India and the Governor General, and the governors of Madras and Bombay, relating to the cultivation of chinchona plants, from April 1866 to April 1870."; House of Commons Parliamentary Papers

1876: East India (chinchona cultivation). Return to an address of the honourable the House of Commons, dated 8 July 1875;—for, copies of the chinchona correspondence (in continuation of return of 1870) from August 1870 to July 1875.; House of Commons Parliamentary Papers

1866: East India (chinchona plant). Return to an address of the Honourable the House of Commons, dated 14 May 1866;—for, "copy of further correspondence relating to the introduction of the chinchona plant into India, and to proceedings connected with its cultivation, from April 1863 to April 1866."; House of Commons Parliamentary Papers

1877: East India (Chinchona cultivation). Further return to an address of the honourable the House of Commons, dated 8 July 1875;--for, copies of the Chinchona correspondence (in continuation of return of 1870) from August 1870 to July 1875.; House of Commons Parliamentary Papers

■ 1887: Ceylon in the "Jubilee year."; Miscellaneous

1883: Papers relating to Her Majesty's colonial possessions. Reports for 1882. (In continuation of [C.-3642.] of June 1883.); House of Commons Parliamentary Papers

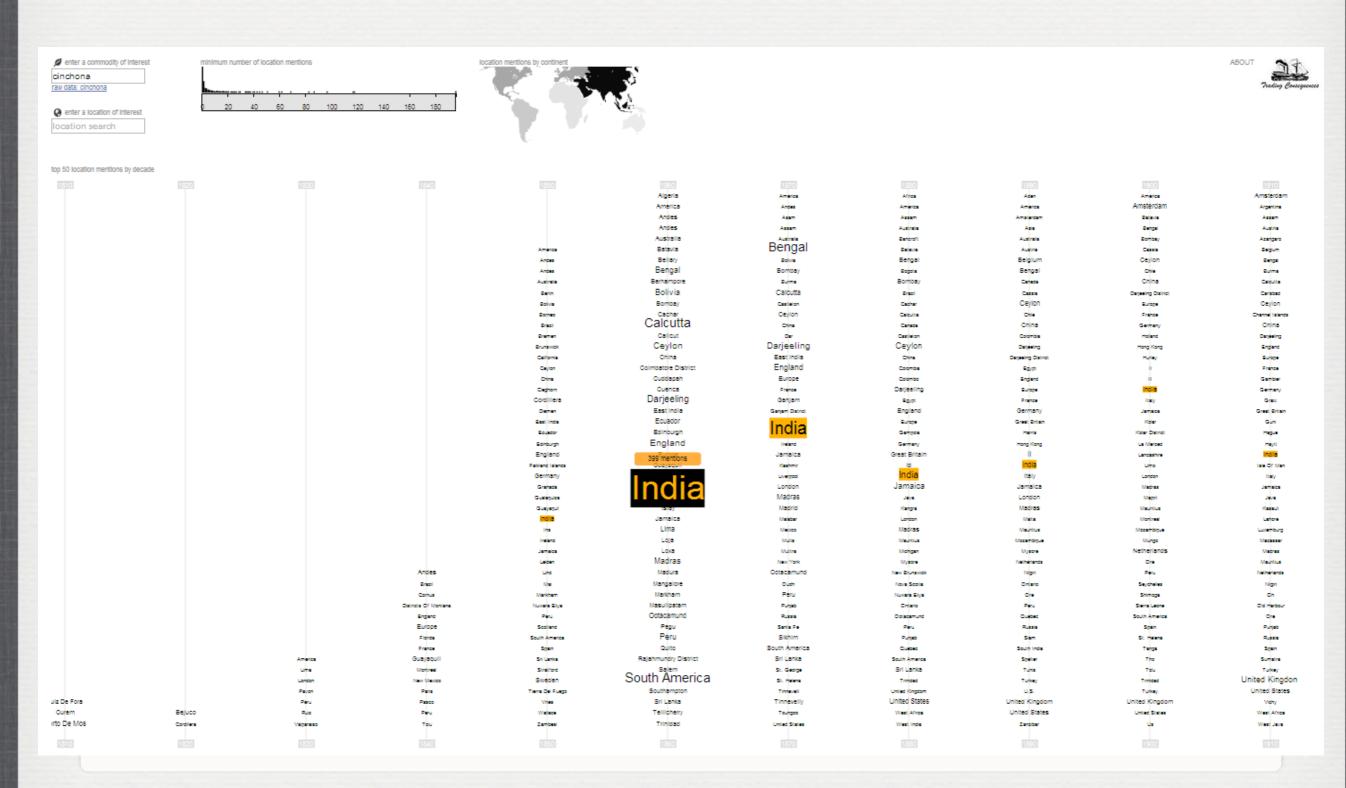
 1885: Ceylon \& Her Planting Enterprize: In Tea, Cacao, Cardamoms, Cinchona, Coconut, and Areca Palms ...; Miscellaneous

■ 1921: Trade and navigation. Return (in part) to an order of the Honourable the House of Commons, dated 16 February 1921;—for accounts relating to trade and navigation of the United Kingdom, for each month during the year 1921.; House of Commons Parliamentary Papers

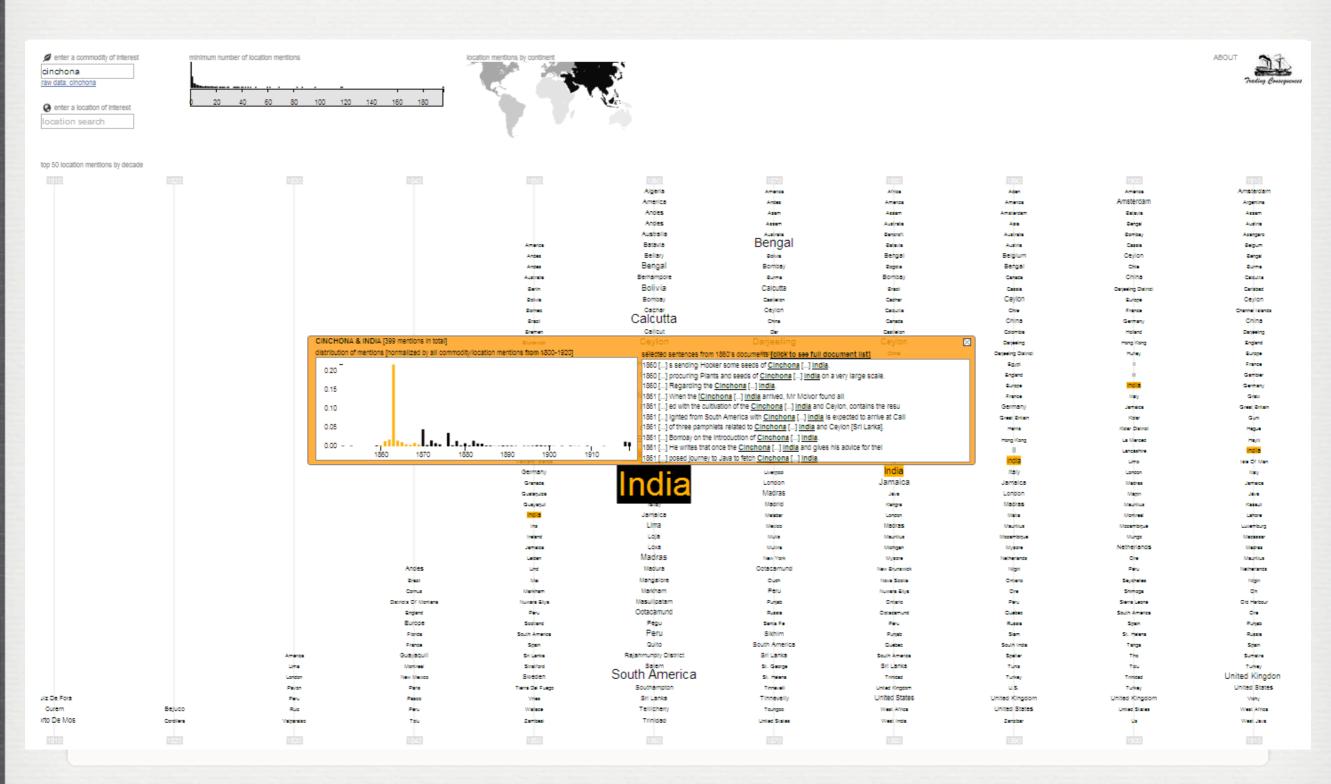
1919: East India (Industrial Commission, 1916 18). Minutes of evidence taken before the Indian Industrial Commission, 1916-18. Vol. I.--Delhi, United Provinces, and Bihar and Orissa.; House of Commons Parlamentary Papers

1861: Notes on the medicinal Cinchona barks of New Granada by H. Karsten; and on the Cinchona trees of Huanuco (in Peru);

BRINGING ARCHIVES ALIVE



BRINGING ARCHIVES ALIVE



Search by:

Commodity

Location

Edinburgh

In Country GB

Feature Type Capital Of Top-Level

Administrative Division

Population 435,791

GeoNames Entry View entry



▼ Filter

BY COLLECTION

ΑII

House of Commons Parliamentary Papers (116)

BY DECADE

ΑII

1850s (116)

BY COMMODITY

Gold (116)

Documents in which 'Edinburgh' is mentioned in relation to commodities (Page 1 of 1)

Filtered by: Decade (1850)

Collection (House of Commons Parliamentary Papers)

Commodity (Gold)

Mentions Document Title

- 90 Report from the Select Committee on the Bank Acts; together with the proceedings of the committee, minutes of evidence, appendix and index.
- Twenty-ninth report of the Commissioners of Her Majesty's Woods, Forests and Land Revenues: in obedience to the acts of 10 George IV. (cap. 50), and 2 William IV. (cap. 1).
- 4 Parliamentary Papers. List of the bills, reports, estimates, and accounts and papers, printed by order of the House of Commons, and of the papers presented by command

SUMMARY

- Scholars potentially have access to enormous amounts of data but cannot always easily manage and navigate it.
- Text mining can be applied to process large text collections, enrich existing text with information or pull out trends which can be visualised. It is a way to enable distant reading, even if such technology is not 100% accurate.
- OCR errors in digitised collections can skew results.
- Interdisciplinary setup of Trading Consequences made it more successful for everyone involved. It wouldn't have been possible without the original data.

WHAT CAN PQ DO?

- Sharing OCRed full text data with mining research initiatives similar to Trading Consequences.
- Improve process for arranging legal agreements for sharing this data.
- Enable a feedback mechanism to improve the OCR and ultimately improve search results.

PALIMPSEST: LITERARY EDINBURGH

- Current AHRC big data project: Exploring place in literature by mining and visualising literature set in Edinburgh, (University of Edinburgh, EDINA, University of St. Andrews).
 - Aiming to retrieve all out-of-copy-right literature set in Edinburgh.
 - Developing a fine-grained gazetteer for Edinburgh to enable geo-referencing on the street and building level.

THANK YOU



- Website: http://tradingconsequences.blogs.edina.ac.uk/
- Demo: http://tcqdev.edina.ac.uk/search/commodity/, http://tcqdev.edina.ac.uk/vis/tradConVis
- Contact: balex@inf.ed.ac.uk