



**History**  
Seminar Podcasts and Online Training

**SPOT**

Digital Tools Modules: Semantic Data and Text Mining

Dr Matthew Phillpott  
Institute of Historical Research  
21 June 2012

## **What is Semantic Data?**

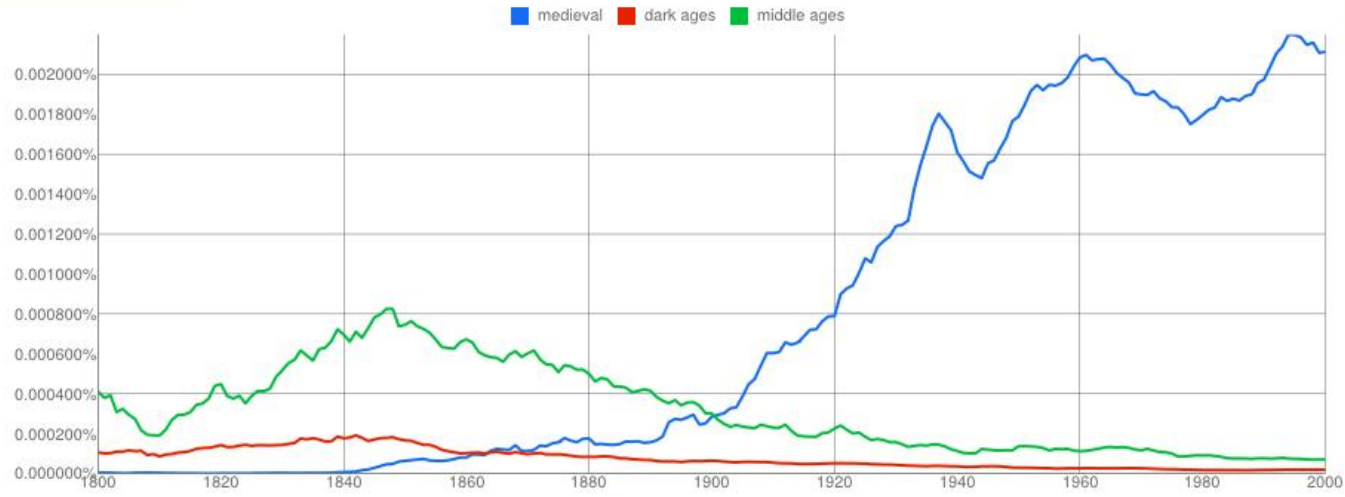
Semantic data is data marked up, however lightly or heavily, in ways which reflect the semantic content of a text, rather than its structure.

## **What is Text Mining?**

An automated method to draw out content based on meaning and context from a large body of text. Enables structure and associations to be revealed in a large body of unstructured data.

# Google books Ngram Viewer

Graph these **case-sensitive** comma-separated phrases:  between  and  from the corpus  with smoothing of .



Search in Google Books:

<a href="#">1800 - 1814</a>	<a href="#">1815 - 1845</a>	<a href="#">1846 - 1854</a>	<a href="#">1855 - 1942</a>	<a href="#">1943 - 2000</a>	<a href="#">dark ages (English)</a>
<a href="#">1800 - 1916</a>	<a href="#">1917 - 1960</a>	<a href="#">1961 - 1968</a>	<a href="#">1969 - 1994</a>	<a href="#">1995 - 2000</a>	<a href="#">medieval (English)</a>
<a href="#">1800 - 1820</a>	<a href="#">1821 - 1844</a>	<a href="#">1845 - 1852</a>	<a href="#">1853 - 1939</a>	<a href="#">1940 - 2000</a>	<a href="#">middle ages (English)</a>

Run your own experiment! Raw data is available for download [here](#).



## Latest podcasts

15 November 2011

[Does the Digital change anything?](#)

15 November 2011

[Great wealth in Argentina, 1810-1930](#)

14 November 2011

[Recreational Music-Making and the Fashioning of Political or Diplomatic...](#)

11 November 2011

[Biology, Brain Theory and History](#)

[more](#)

## New Podcasts: Does Digital Change anything?



**Valerie Johnson and David Thomas discuss**

## News and Updates from History SPOT

### News & Updates

IHRDigProjects Does Digital change anything?  
[#News](https://t.co/88A03FQ9)  
55 minutes ago · [reply](#) · [retweet](#) · [favorite](#)

IHRDigProjects New podcast on wealth in Argentina between 1810-1930: [#News](https://t.co/VxPbdFZ9)  
4 days ago · [reply](#) · [retweet](#) · [favorite](#)

IHRDigProjects Music-making at the court of Elizabeth I - new podcast from History SPOT:  
[#News](https://t.co/XreHjGgI)  
5 days ago · [reply](#) · [retweet](#) · [favorite](#)

IHRDigProjects What does Biology, Brain Theory and



[Join the conversation](#)

## Welcome to History SPOT (Seminar Podcasts and Online Training)

[Podcasts](#) [Research handbooks](#) [Research training](#) [Collaboration](#)

History SPOT presents podcasts from the internationally renowned research seminar programmes hosted by the Institute of Historical Research (IHR) in an interactive and collaborative space. The 'Seminars' section of the site holds our searchable archive of these podcasts from 2009 to the present, alongside discussion forums and other additional content. This section of the site is designed not simply to provide access to this archive but to allow you to comment on and interact with it. If you have a question or would like to start up or enter a debate or discussion about the topic you can do so here with other historians.

- Archived podcasts from 2009 to the present selected from over 50 history research seminar programmes held at the



## Module Index

1. Introduction
2. Creating and gathering data
3. Regular Expressions and Scripting
4. Text Analysis Techniques
5. Named Entry Recognition
6. Taking things Further
7. Further Reading


## Digital Tools Home

[Return to Digital Tools home page](#)

## Settings

- < Course administration
  - > [Question bank](#)

---

- < Switch role to...
  -  [Return to my normal role](#)

---

- > [My profile settings](#)

---

- > [Site administration](#)

## Topic outline

### Text Mining for Historians

#### 3 Regular Expressions and Scripting

Duration: 1.5 hours

Author: [Dr Mark Merry](#)

#### [Contents and Learning Outcomes](#) (click here)

##### 1.

This section of the course introduces you to regular expressions. Understanding these is essential to being able to use electronic texts for text mining. Reading through the handbook and the exercises that are contained within.

#### [Regular Expressions and Scripting](#)



[Previous](#)



[Next](#)

Jump to...

## Table Of Contents

1. Why mark up text?
2. What are markup languages?
3. HTML
  - 3.1 Marking up HTML
  - 3.2 Exercise I: HTML mark-up
  - 3.3 Entity references
  - 3.4 Exercise II: HTML mark-up
  - 3.5 Exercise III: Parish Clerk's Memorandum Book
4. Further resources



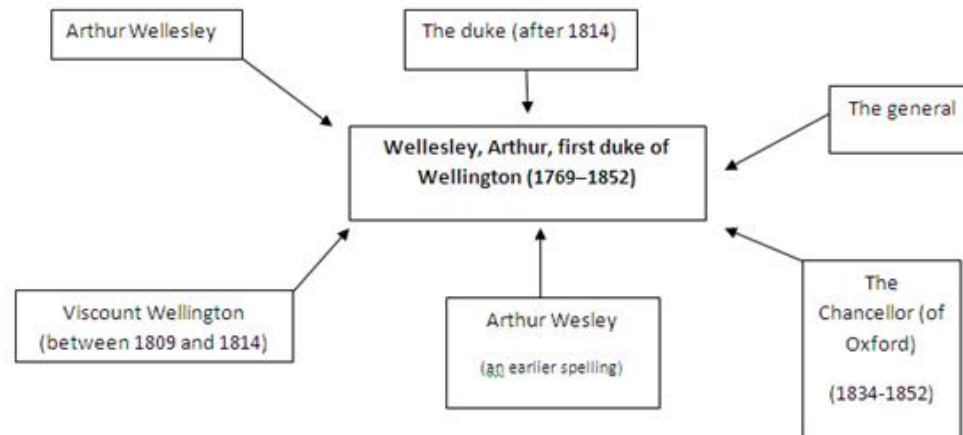
## 1. Why mark up text?

In a file without markup, the running text will normally be in **plain text** format. To take an extremely simple example, imagine a **text file** that just contains

*In Wellington I saw a statue of Wellington*

If we want to search this kind of text file for the string Wellington we can easily find both of the examples by using a find command. But there is no way to distinguish between the person and the place.

Suppose that we want to find all references to the First Duke of Wellington in a text (one that contains, rather than the eight words above, a couple of million words). We don't want references to places called Wellington, or to the duke's son, the Second Duke of Wellington, or any other people called Wellington. But we do want to find all of these references to the same man:



Marking up a text semantically is the most reliable way to be able to return all of the information you want to get from your search. People, places and dates are often the focus of semantic markup, but anything can be marked up - quotations, emotions, economic data or food - anything that is of interest to the researcher. In this course we'll show you some ways of doing this.

### Module Index

- 1. Introduction
- 2. XML I
- 3. XML 2
- 4. Markup Schemes
- 5. Project advice

### Digital Tools Home

Return to Digital Tools home page

### Settings

- > Course administration
    - > Question bank
- > Switch role to...
    - Return to my normal role
- > My profile settings
- > Site administration



## Topic outline

### Semantic Data for Historians

#### 2 XML I



Duration: 2 hours

Author: [Jonathan Blaney](#)

#### Contents and Learning Outcomes [\(click here\)](#)

##### 1.

First listen to this video clip taken from one of the research seminars held at the IHR. In this clip Mangus Huber talks briefly about semantic markup for the Old Bailey criminal trials for 18th and 19th century London.

The video player shows a snippet of XML code from 'Original computerized Proceedings (Sheffield)'. The code defines a person named Sarah Sanders, indicted for stealing a Portugal Piece of Gold, value 36 s, a Gold Ring, value 10 s, a Gold Ring set with Vermillion Stones, value 7 s. 6d, a Silver Girdle Buckle, value 10 s, three Aprons, a Shirt, a Shift, and 2 Ells of Holland, the Goods of John Underwood, in his House, on March 4.

The video player interface includes a play button, a progress bar at 01:31, and the Vimeo logo.

To listen to the full video or audio podcast click here



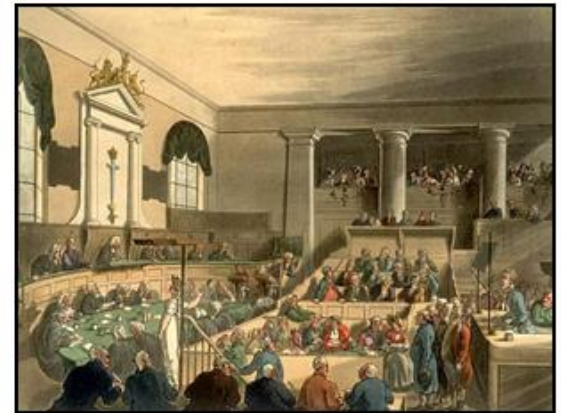
## Table Of Contents


1. What is Text Mining?
  - 1.1 Example: The Old Bailey Proceedings
  - 1.2 Uses of Text Mining
  - 1.3 How does text mining work?
  - 1.4 What are the limitations and barriers of text mining?
2. Case Study: Ngram Viewer
  - 2.1 What does the Google ngram viewer do?
  - 2.2 Example: "John Foxe"
  - 2.3 Language options
  - 2.4 Exercise: Try it yourself

## 1. What is Text Mining?

### 1.1 Example: The Old Bailey Proceedings

For example, would you like to know how often the word 'guilty' appears in the Old Bailey trial accounts? The answer is findable using a standard search engine on the Old Bailey Online website (the answer is **182612**). How about how many people were found guilty? The answer is **163261**. What about the number of defendants found guilty of murder? The answer is 1518. These last two figures are not possible to find through the standard search engine as they are an entirely different type of question; we are not looking for how many times the word 'guilty' appears in the proceedings but how many trials resulted in a guilty verdict. We want to discover something meaningful within the body of texts automatically rather than manually checking each and every trial account.



This is a relatively simple example of data mining where the original documents have been marked up and tagged by surname, given name, alias, offence, verdict, and punishment. To calculate those results manually you would have to work your way through 197,745 criminal trial accounts (some 127 million words in total). 

This form of text mining, however, is little more than an advanced search engine. Useful, yes but limited. As the creators of the Online Old Bailey themselves admit (and have attempted to redress in a subsequent project):

'Analyzing this kind of data by decade, or trial type, or defendant gender etc., can re-enforce the categories, the assumptions, and the prejudices the user brings to each search and those applied by the team that provided the XML markup when the digital archive was first created'.



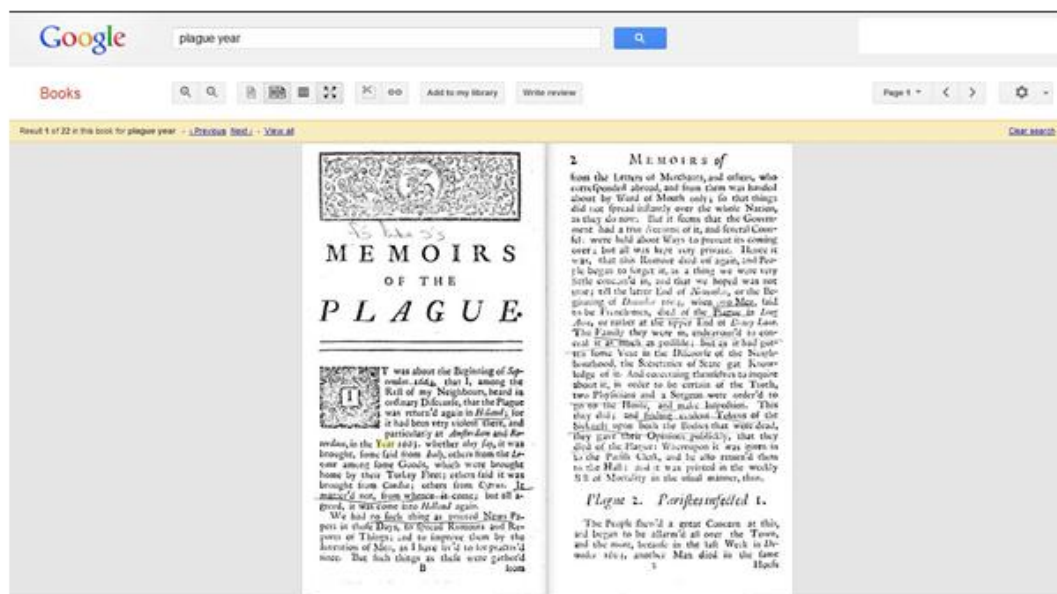
## Table Of Contents

1. Downloading Data
  - 1.1 Locating books on Google Books
  - 1.2 Downloading books in Google Books
  - 1.3 Exercise: Google Books
  - 1.4 Other downloadable books online
2. Using data from multiple sources
  - 2.1 Project Gutenberg
  - 2.2 Oxford Text Archive
3. Creating and re-using data
  - 3.1 Optical Character Recognition (OCR)
  - 3.2 Exercise: OCR
  - 3.3 Rekeying and crowdsourcing

## 1. Downloading Data

### 1.1 Locating books on Google Books

If we continue with the example of the texts available through Google Books, we might try to use the text location tools on offer to find books about the devastating outbreak of plague in London in 1665. Naturally enough we would end up sooner or later at Defoe's (largely fictional) account *A Journal of the Plague Year*.



Daniel Defoe's *A Journal of the Plague Year*, 1722 (Google Books)

When viewing a book online in the Google Books collection, assuming that a digital copy is available (and not just a catalogue record, which is sometimes the case) you are presented with a scanned 'facsimile' of the book. This can often tell you quite a lot about the quality of the original copy of the book that has been used in the digitisation – note in this case the scribbles and underlining of several passages. Of particular interest is the warning that appears above the title from the well-meaning owner of the scanned copy, an issue we will be

Donna Baillie talking about her experience on ReScript:

[Audio File](#)

Dan Cohen talking about the Ngram Viewer:

[Video File](#)

History  
Seminar Podcasts and Online Training

SPOT



Semantic Data and Text Mining for Historians Training Modules:

<http://historyspot.org.uk>

Coming Soon