



#LancsBox X

Corpus building

Exercises: building a corpus & exploring the results

In this handout, we will practice two methods of building a corpus: a manual method and an automatic method using #LancsBox X. You will then explore the data in #LancsBox X using the Text tool. The aim of this worksheet is for you to learn

- to understand the process of corpus creation
- to evaluate the quality of the data
- to understand key parameters of texts in a corpus

Text

corpus creation

Web

tokens

Task 1

Build your own newspaper mini-corpus

- a) Decide on the newspaper you want to use. Choose ONE of the following British newspapers:

Newspaper	Website
The Guardian	www.theguardian.com
The Daily Mail	www.dailymail.co.uk

- b) Choose a word or phrase that characterizes a topic of your interest (e.g. “English teaching”, “ecology”, “school tests”, “healthy lifestyle”, “poverty”, “climate change”, etc.).

My word/phrase is: _____

- c) Go to [Google](#) or another search engine to search for articles in the selected newspaper. Include your word/phrase(s) in the search,

e.g. "climate change" site:www.theguardian.com

NB: There is no space between *site* and the web address.

- d) Open the articles returned by Google (or other search engine) one-by-one and copy-paste each into a separate text document. Do this with at least 10 articles.
- e) Save your text documents as plain text (.txt) in a folder. You can use your own text editor or download the free [Sublime Text](#). MS Word or similar word processors are not ideal but will do the job for this exercise; just remember to save the files as plain text choosing the encoding “Unicode (UTF-8)” in the saving options.

Google tip: If you get many results, you can limit your search to a particular period of time by clicking on “Tools” and then on “Custom range”

Customised date range

From « December 2020 »

To

M	T	W	T	F	S	S
30	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	1	2	3
4	5	6	7	8	9	10

Congratulations! You’ve just created your own mini-corpus!

Task 2

Uploading your corpus into #LancsBox X

- a) Open #LancsBox X and make sure that all your files are in a single folder e.g. MY_CORPUS. The data does not need to be in any special format – the tool accepts all major data formats (e.g. .txt, .doc., .pdf, etc.). From the corpus drop down menu choose ‘add corpora’ and select ‘My data’.

- b) Explore the corpus in the Text tool . Summary statistics are available when you hover with your mouse over the name of the corpus.

```

102,305,246 grammar tokens
99,949,544 space tokens
88,171 texts
Full name: The British National Corpus 2014
Version: 2.0 CLAWS7
Annotations: CLAWS7 & USAS
Folder: BNC2014_2
    
```

Fill in the table below with the descriptive statistics of your mini-corpus:

Corpus size – space tokens (strings separated by space)	
Corpus size – grammar tokens (identified by the tagger)	
Texts	

- c) Tagging. When you load your own corpus into #LancsBox X, you can select ‘Grammatical’ and ‘Semantic’ tagging, which is done automatically. In the Text tool, double-click on any text and view all available annotation

Task 3

Exploring your corpus in #LancsBox X

- a) Go to the Text tool in #LancsBox X and search for the frequency of your word or phrase from Task 1 in your mini-corpus.

The screenshot shows the 'Text' tool interface. At the top, there's a search bar with 'web' selected and a dropdown menu set to 'whole corpus' with '280K' next to it. Below this, a search bar contains 'research'. To the right of the search bar, it displays 'Hits: 216 (771.41)' and 'Texts: 47/100'. A callout box labeled 'Summary information' points to these statistics. Below the search bar, there's an 'Overview' section showing '100' items. A table below that lists search results with columns: Name, Tokens, Frequency, Rel. freq..., and MATTR₅₀. A callout box labeled 'Individual text info' points to the first row of the table.

Name	Tokens	Frequency	Rel. freq...	MATTR ₅₀
https://en.wikipedia.o...	931	27	29,001.07	

How many times does your selected word/phrase occur in the corpus? _____

In how many texts does it appear? _____

What is the text with the highest relative frequency of your word or phrase? _____

Task 4

Creating a corpus automatically in #LancsBox

#LancsBox X offers a feature, which allows automatic download of texts from the web. In this task, you will explore this feature. Imagine that you want to download a corpus made of Wikipedia articles related to linguistics.

- a) Locate your initial (seed) website that includes the landing page about linguistics: <https://en.wikipedia.org/wiki/Linguistics> From the corpus drop down menu choose 'add corpora' and select 'Web'.

The screenshot shows the 'Add corpora' form with three tabs: 'Corpus hub', 'My data', and 'Web'. The 'Web' tab is active. It contains fields for 'Corpus full name*' (with 'Wiki_linguistics' entered), 'Short display name', 'Language' (set to 'English'), and 'Initial URLs*' (with 'https://en.wikipedia.org/wiki/Linguistics' entered). There are 'Create corpus' and 'Close' buttons at the bottom. Three callout boxes provide instructions: '1. Name your corpus' points to the 'Corpus full name' field, '2. Paste the Wikipedia URL' points to the 'Initial URLs' field, and '3. Click on "Create corpus"' points to the 'Create corpus' button.

- b) Choose a different corpus name e.g. 'Wiki_linguistics_tagged'. Repeat the process with Grammatical and Semantic tagging on.

Tagging Grammatical Semantic

Select the appropriate tagging. If you are using this feature for the first time, the appropriate taggers and language models will need to be downloaded; this is done automatically.

Task 5

Exploring your Wikipedia corpus in #LancsBox

- a) Explore the size of your corpus (see instructions in Task 2b). Fill in the table below with the descriptive statistics of your Wikipedia corpus:

Corpus size – space tokens (strings separated by space)	
Corpus size – grammar tokens (identified by the tagger)	
Texts	

- b) How many times do different words related to linguistics occur in the corpus?

Search term	Absolute (relative) frequency	No. of texts
language		
corpus linguistics		
...		
...		



Developed at
**Lancaster
University**

