Lancaster
University

# Doing corpus linguistics with #LancsBox

In these tasks, we will become familiar with some of the functions of #LancsBox by investigating features of modern-day British English using the BNC2014 Baby+.

**Tasks 1 - 2** will focus on practicing different types of **searches** using the KWIC tool.

**Tasks 3 - 5** will introduce how to create and change the settings of **collocation graphs** using the GraphColl tool.

**Tasks 6 - 7** will introduce how to **create a web corpus automatically** (a brand-new feature for #LancsBox v.6).

There are also **optional tasks** you can try during or after the practical session.

We will be using the **British National Corpus 2014 Baby+ edition;** this corpus is a four-million-word balanced subset of the BNC2014. Further details can be found here: http://corpora.lancs.ac.uk/lancsbox/docs/pdf/BNC2014Baby.pdf

**Task 1. Searches.** Go to the KWIC tool in #LancsBox and search for the following expressions in the BNC2014 Baby (provided with #LancsBox). Note down their frequencies and distributions in texts.

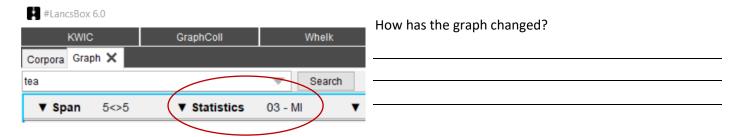| Type of search | Search term | Occurrences (per 10k) | Number of texts |
|---|---|---|---|
| Simple | tea | | |
| Simple | weather | | |
| Phrase | mug of tea | | |
| Wildcard | rain* | | |
| Smart Search | DOWNTONER | | |
| Regex | /mate\|friend/ | | |
| Regex | nick [as headword] V* [as POS] | | |

**Optional task:** You are researching how 'hailnames' (informal forms of address) are used in British English. How might you build one query to search for *pal*, *mate* and *buddy* simultaneously? What are some issues you might face?

**Task 2. Applying filters.** Still in the KWIC tool, search for the following expressions and apply filters. Note down their frequencies and distribution in texts.

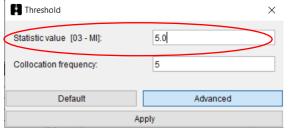| Search term | Filter | Occurrences (per 10k) | Number of texts |
|---|---|---|---|
| tea | make [anywhere LEFT] | | |
| fish | chips [in R2 position] | | |

**Task 3. Create a collocation graph and change settings.** Go to the GraphColl tool, follow the directions and note down the frequencies.

(a) Build a collocation graph by conducting a simple search for *tea.* What results did you get?

(b) Change the statistical association measure to MI score. This will delete the current graph. Search for *tea* again to create a new graph.
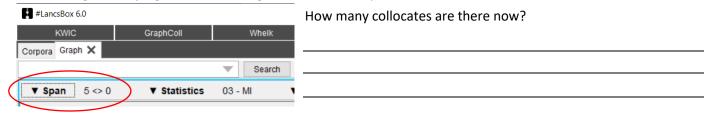
#LancsBox v.6 manual: http://corpora.lancs.ac.uk/lancsbox/docs/pdf/LancsBox_6.0_manual.pdf

How has the graph changed?

_____

_____

_____

(c) Open the <u>threshold settings</u> and change the statistic value to 5.0 (3 and higher is considered strong for the MI score equation). Search for _tea_ again.



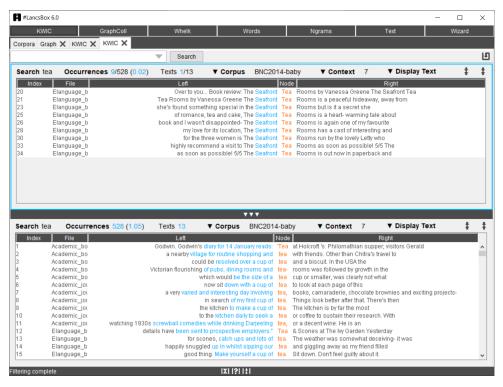How has the graph changed now?

_____

_____

_____

(d) Change the <u>window span</u> to search for five words to the left and zero words to the right (5L, 0R). Search for _tea_ again, keeping the same settings from the last steps.



How many collocates are there now?

_____

_____

_____

**Optional task:** You are researching what collocations British English speakers use when talking about the weather. How might you start this study using GraphColl? What settings would you consider changing, and why? Why might you use a narrower or wider collocation window?

**Task 4. Combining GraphColl and KWIC view to explore collocation context**. To help interpret graphs, GraphColl offers a concordance function, which displays examples of collocate use (KWIC). To display the concordance lines of a given collocation, **right click** on the collocate in the table or in the graph. These concordance lines can be expanded into the full-screen KWIC view by clicking on three arrows ( ▲▲▲ ) at the top right. In the top panel, the full-screen view displays examples of the selected collocate co-occurring with the node; in the bottom panel, all collocates are displayed.

#LancsBox v.6 manual: http://corpora.lancs.ac.uk/lancsbox/docs/pdf/LancsBox_6.0_manual.pdf

Keeping the same settings, search for *tea* and right click on the collocate *seafront*. Explore the context in which it occurs. Comment on the association between *tea* and *seafront* in this corpus.

_____

_____

_____

**Task 5. Build collocation networks.** GraphColl also offers the function to explore second-order collocates through building collocation networks. These are visualisations that help you explore language patterns.

(a) Search for *tea* to create a new graph (use MI score and settings: span 5L, 5R; Statistic value = 5.0; Collocation frequency = 15).

(b) Find the collocate *cup* in the graph (or table) and double click on it. This should create a collocation network similar to the one below. (n.b. you can zoom into a graph using your mouse wheel or change the text size using Ctrl +/-)

**|<---------------------------------------second-order collocation network-------------------------------------------------->|**



**|<------first-order collocation network------>|**

**|<---------------shared collocates--------------->|**

#LancsBox v.6 manual: http://corpora.lancs.ac.uk/lancsbox/docs/pdf/LancsBox_6.0_manual.pdf

(c) Comment on the connection between the collocates. Feel free to explore the network further by clicking on any of the collocates.

_____

_____

_____

**Optional task:** You want to investigate both collocates *cup* and *cups*. To do this, change the unit setting from type to lemma.

| ▼ Span | 5<>5 | ▼ Statistics | 03 - MI | ▼ Threshold | ▼ Corpus | BNC2014-baby | ▼ Lemma |

How does the graph change? Explore further by changing the span, statistics or threshold settings as you like.

_____

_____

_____

**Task 6. Create a web corpus automatically.** #LancsBox v.6 offers a new feature, which allows automatic download of texts from the web.

(a) Locate your initial (seed) website that includes the landing page about your chosen topic.
(b) On 'Corpora' tab, click on 'Corpus' under create (see image)
(c) Paste the URL into the URL box and click on 'Create corpus'
(d) Wait for the process to finish – by default #LancsBox is downloading and converting into text 100 websites at 2 levels of embedding – and explore the log and the folder with the data. Click 'OK'.
(e) Click on 'Import' to import and automatically POS-tag your corpus.

**Task 7. Explore and search your corpus using #LancsBox.**

(a) Explore the size of your corpus and note it down:

| | |
|---|---|
| **Tokens (running words):** | |
| **Types (different words):** | |
| **Lemmas (headword + POS category):** | |

(b) Explore your corpus using relevant search terms:

| Search term | Occurrences (per 10k) | Number of texts |
|---|---|---|
| | | |
| | | |
| | | |
| | | |