# #LancsBox X

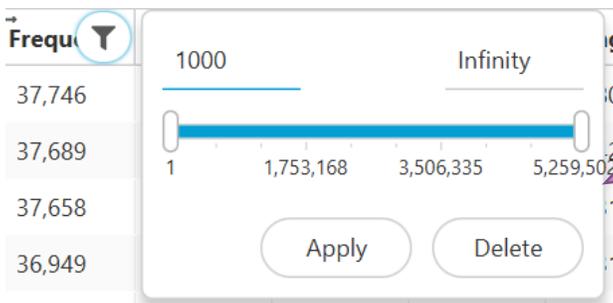# Words

In this task, you will explore the information about the frequency of words. Go to the Words tool in #LancsBox X and select the BNC2014 corpus (whole corpus). Keep the unit as 'word (lowercase)' and find the following information:

1. Number of words in the wordlist _____858,450_____

2. The most frequent word in the wordlist is _____the_____

3. The information you can find about the most frequent word__Absolute Freq., Relative Freq., Dispersion__

4. Number of words with a frequency ≥1000 _____7,333_____

5. Number of words with a frequency ≥100 _____32,672_____

6. Number of words with a frequency ≥10 _____120,571_____

7. Number of words with a frequency of 1 (so called *hapax legomena*) _____488,054_____

| Freque | | g |
|--------|--|---|
| 37,746 | 1000      Infinity | |
| 37,689 | | |
| 37,658 | 1    1,753,168   3,506,335   5,259,50 | |
| 36,949 | Apply      Delete | |

**Tip**: Click on the filter icon and apply an appropriate frequency range to the 'Frequency' column. Click 'Apply'.

## Words: Terminology

There are different concepts of a word. In corpus linguistics, terms such as *token, type, lemma* or *lexeme* are often used to denote different senses in which the general term 'word' is used.

**Token (running word)** is a single occurrence of a word form in a text or corpus.

**Type** is a unique word form in a corpus.

**Lemma** denotes all inflected forms belonging to one stem and one word class; in #LancsBox by default, a combination of a headword and a grammatical category (e.g. go + VERB). For example, a lemma 'go' includes the following word forms (types): 'go', 'goes', 'went', 'going' and 'gone'.

**Lexeme** is a lemma with a particular meaning attached to it, which is necessary to distinguish polysemous words (words with multiple meanings).

In this task, you will explore the information about the distribution of words in texts. This is what is called 'dispersion' in corpus linguistics. You will be using BNC2014 (whole corpus).

**2a**    What words do you think appear in most English texts? Give some examples.

_____

_____

**2b**    Based on your intuition, try to guess approximately the percentage of texts the following words occur in:

*research _____% cake _____% think_____% hitherto_____% February _____%*

### Now check your answers using corpus data:

**2c**    In Words, sort the table according to 'Range %' and note down top five words and the percentage of texts they appear in.

The (99.15) , to (97.05) , a (96.78) , and (96.57) , of (95.85)

_____

**2d**    Now search for the following words and note down the percentage of texts they occur in:

*research 8.86%    cake _1.74%    think_22.65%    hitherto_0.31%    february _5.11%*

**2e**    Now search for the following words and note down their DP values:

*research DP_0.77_ cake DP__0.88_    think DP_0.52_    hitherto DP_0.98_    february DP_0.88__*

**Tip**: DP stands for 'Deviation of proportions'. It operates on a scale 0 – 1 with 0 being the most equally distributed an 1 being extremely unequally distributed.

| DP (deviation of ... |
| --- |
| 0.15 |
| 0.11 |

In this task, you will explore different units that a wordlist can be based on. By default, #LancsBox X shows you lowercase types 'word (lowercase)'. To obtain even more information, the units can be modified.

**3a** In turn, change to unit to 'class', 'pos' and 'sem' and fill in the table below to identify top 5 wordclasses, pos-tags and semantic tags.

| Top 5 word classes | Top 5 pos tags | Top 5 semantic tags |
|---|---|---|
| SUBST | NN1 | Z5 |
| ADJ | JJ | Z8 |
| VERB | II | Z99 |
| PREP | AT | A3 |
| ART | NN2 | Z1 |

**Tip**: In Words, change the unit in the top bar next to the corpus information.
For the interpretation of pos tags go to: https://ucrel.lancs.ac.uk/claws7tags.html
For the interpretation of semantic tags or to: https://ucrel.lancs.ac.uk/usas/USASSemanticTagset.pdf

W°
BNC2014 2.0 CLAWS7     whole corpus     ▼     102M     class     ▼

**3b** Now, change the unit to lemma. In turn, apply filters to the 'Term' column to identify most frequent nouns _N, verbs V_, and adjectives _J.

| Top 5 nouns _N | Top 5 verbs _V | Top 5 adjectives _J |
|---|---|---|
| year | be | good |
| time | have | new |
| people | do | other |
| thing | say | great |
| way | will | high |

**3c** Interpret the findings from 3b 'Top 5 nouns'. What are the concepts that are most frequently discussed in current British English?

Time, people and general concepts – these are typically useful across a range of contexts.
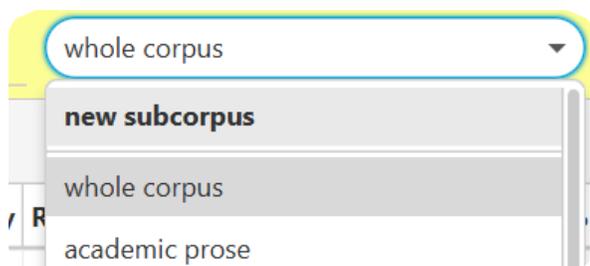
It is interesting to see that similar terms - certainly from these areas - will occur both in speech and writing, reflecting the need to frame discourse in terns of time, space, relevant actors (people; down the list: man, woman), and general objects and processes (thing, way).

Keywords are words that occur with a considerably higher frequency in a given (sub)corpus compared to a reference (sub)corpus. In this task, you will explore keywords used in SMS messages compared to all of elanguage.

1. First create a SMS messages subcorpus of the BNC2014 and note down its size: ___228k_____

2. Click on the keyword icon 🔑 and select **BNC2014 elanguage** as a reference corpus.

3. Note down the top 10 keywords, i.e. words typical of SMS messages _____u, lol, haha, yeah, yep, ok, yeh, tho, av, oh _____

4. Note especially the 1st keyword _____u_____ and its meaning ____2nd person pronoun 'you'

5. Note also the relative frequency of the 1st keyword in SMS messages _____4727.86_____ compared to its frequency in the reference corpus (whole elanguage) ___302.56_____.

| whole corpus ▼ |
|---|
| **new subcorpus** |
| whole corpus |
| academic prose |

**Tip**: To create a subcorpus, go to the subcorpus manu and click on 'new subcorpus'. Then select the appropriate category (e.g. genre>elanguage, subgenre> SMS messages)

Developed at
Lancaster University