# THE AMC-50 DATA QUICK REFERENCE

The **DATA REFERENCE** contains MDA DATA from previous studies - see *Opening Credits: Multi-Dimensional Analysis (MDA) DATA* - which uncover the linguistic and textual features of movie discourse and its similarity with face-to-face conversation and NEW DATA from the dialogs of 50 movies (henceforth AMC-50 DATA) produced in the USA from 1959 to 2019 (see Forchini 2021:33-35 and Table 1 here). The AMC-50 DATA have been processed and extracted via AntConc 4.0.5 (www.laurenceanthony.net/software/antconc), #LancsBox 6.0 (http://corpora.lancs.ac.uk/lancsbox), SketchEngine (www.sketchengine.eu), *WordSmith Tools* 8.0 (https://lexically.net) and *RStudio* (2020).

Table 1. The AMC-50 movies

## THE AMC-50 SIZE

The size of the AMC varies depending on what is being counted and what counts as a word: a TOKEN is the smallest unit that a corpus consists of, consequently, the number of tokens in a corpus represents the total number of individual words it contains. A TYPE, instead, is a unique word form in a corpus, consequently, the number of types in a corpus represents the number of unique word forms it contains.

| SOFTWARE | TOKENS | TYPES |
|---|---|---|
| AntConc 4.0.5 | ≈ 560000 | 18216 |
| #LancsBox 6.0 | ≈ 530000 | 20477 |
| SketchEngine | ≈ 580000 | 19950 |

# CLAPPERBOARD

**LEGAL DISCLAIMER:**
The **A**merican **M**ovie **C**orpus (AMC) is conceived as a repository of data transformed into a new monomodal (textual) utility. As such, the AMC is a collection of movie dialogs transcribed by the AMC team which does not include any audiovisual (multimodal) material or script from the web. The copyright of the movies resides with the original copyright holder(s). The resulting texts are used for noncommercial/nonprofit purposes, such as linguistic research, scholarship, teaching, criticism and comment. The AMC is not publicly available, but data extracted from the corpus can be shared free of charge. Any scholars, language learners and/or teachers (henceforth USERS) who are interested in word lists, lexical bundles and collocations can use the data presented here by citing us. USERS can also access *The AMC LAB* for further updates and contact *The AMC Team* to obtain concordances and snippets of dialog (cf. www.americanmoviecorpus.net). USERS declare and accept that the data obtained from the corpus will be used by them for the exclusive purposes of research, scholarship, teaching, criticism and comment. USERS assume complete responsibility: neither the AMC team nor the AMC Board are a party to or are in any way responsible for any copyright infringement.

**BACKSTORY: The AMC LAB** has been created to share data with scholars, teachers and learners who aim to investigate, teach and/or learn the lexico-grammatical features characterizing spoken language. Although movies are artifacts by nature, recent investigations of the AMC (see Publications on the AMC site) have, in fact, revealed that their dialogs share the same textuality and linguistic features of natural face-to-face conversation. These revolutionary findings have opened up new avenues:

**For SCHOLARLY RESEARCH:** for many years movie language has been considered as artificially written-to-be-spoken and deemed unlikely to comprise the features that characterize conversation. Data from the AMC corpus offers new ways of approaching the study of movie language;

**For LANGUAGE LEARNING:** language learners can improve their spoken competence through practice on movie conversation;

**For LANGUAGE TEACHING:** authoritative scholars have been emphasizing the crucial role played by spoken language in communication for almost a hundred years. In spite of this, attention given to the study of lexico-grammatical spoken features in educational settings has been scarce. The textual and linguistic similarity of movie dialogs with face-to-face conversation and the rich resource of spoken language features which the AMC represents mean that teachers now have the chance to give spoken language its rightful place.

**ASIDE:** The data shared here are intended as a mere QUANTITATIVE REFERENCE for learners, teachers and scholars interested in movie discourse and are just an example of the role that movie language and corpora can have in the mastering of the most recurrent linguistic patterns found in conversation.

## WORD LISTS & MULTI-WORD SEQUENCES

The concept of multi-word sequences, also called lexical items (Sinclair 1998, 2004), clusters (Scott 1998, Scott and Tribble 2006), lexical bundles (Biber et al. 1999), n-grams (www.laurenceanthony.net), sequences of words (Hunston 2006), or phrasal units (Stubbs 2006), is directly linked to the Firthian notion of collocation in that it expands the category in terms of number of words involved: words do not only come in sets of two (as collocates do), but also in groups of three, four (or more) items which together create a meaning which is different from the meaning of the single items taken in isolation (Sinclair 2004). More recent studies have distinguished *lexical bundles* from *n-grams*, defining the former specifically as uninterrupted strings of three or more words which, in order to be considered as such, have to be extremely frequent in a given register (Cortes 2015).

| RANK | AntConc 4.0.5 | | #LancsBox 6.0 | | SketchEngine | |
|---|---|---|---|---|---|---|
| | TYPE | FREQUENCY | TYPE | FREQUENCY | TYPE | FREQUENCY |
| 1 | you | 24214 | you | 21151 | you | 24210 |
| 2 | i | 23336 | i | 16909 | i | 23334 |
| 3 | the | 15149 | the | 15143 | the | 15145 |
| 4 | a | 11579 | a | 11021 | a | 11028 |
| 5 | a | 11074 | to | 10987 | to | 10996 |
| 6 | to | 11013 | and | 7813 | it | 10629 |
| 7 | it | 10611 | it | 7131 | 's | 10230 |
| 8 | that | 8820 | that | 6779 | that | 8820 |
| 9 | and | 7831 | of | 6144 | and | 7815 |
| 10 | t | 7204 | is | 5334 | n't | 7164 |
| 11 | of | 6150 | n't | 5140 | do | 6629 |
| 12 | what | 5632 | me | 4908 | of | 6344 |
| 13 | is | 5332 | what | 4871 | what | 5627 |
| 14 | we | 5240 | this | 4804 | is | 5626 |
| 15 | in | 5167 | no | 4361 | we | 5240 |
| 16 | me | 4915 | oh | 4265 | in | 5146 |
| 17 | this | 4815 | on | 4213 | me | 4912 |
| 18 | in | 4546 | i'm | 3814 | this | 4815 |
| 19 | no | 4372 | your | 3806 | 'm | 4448 |
| 20 | oh | 4271 | we | 3779 | no | 4361 |

| RANK | AntConc 4.0.5 | | #LancsBox 6.0 | | SketchEngine | |
|---|---|---|---|---|---|---|
| | TYPE | FREQ. | TYPE (with space separators) | FREQ. | TYPE | FREQ. | TYPE | FREQ. |
| 1 | i m | 4472 | you know | 1497 | you know | 1497 | do n't | 3222 |
| 2 | it s | 3381 | are you | 1290 | are you | 1288 | i do | 1563 |
| 3 | don t | 3223 | i don | 1279 | in the | 1219 | you know | 1497 |
| 4 | you re | 2... | | | | | |
| 5 | that s | 2... | | | | | |
| 6 | you know | 1... | | | | | |
| 7 | are you | 1... | | | | | |
| 8 | i don | 1... | | | | | |
| 9 | in the | 1... | | | | | |
| 10 | do you | 1... | | | | | |

| RANK | AntConc 4.0.5 | | #LancsBox 6.0 | | SketchEngine | |
|---|---|---|---|---|---|---|
| | TYPE | FREQ. | TYPE (with space separators) | FREQ. | TYPE | FREQ. | TYPE | FREQ. |
| 1 | i don t | 1279 | no no no | 471 | no no no | 486 | i do n't | 1277 |
| 2 | don t know | 561 | what are you | 378 | what are you | 389 | do n't know | 560 |
| 3 | i m not | 519 | oh my god | 309 | i don't know | 372 | you do n't | 489 |
| 4 | you don t | 499 | what do you | 283 | oh my god | 309 | no no no | 477 |

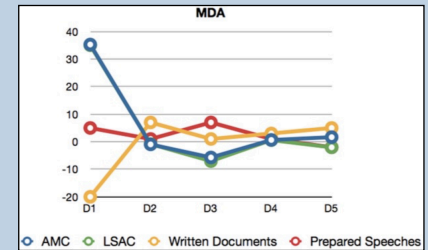| RANK | AntConc 4.0.5 | | #LancsBox 6.0 | | SketchEngine | |
|---|---|---|---|---|---|---|
| | TYPE | FREQ. | TYPE (with space separators) | FREQ. | TYPE | FREQ. | TYPE | FREQ. |
| 1 | i don t know | 441 | no no no no | 252 | no no no no | 257 | i do n't know | 440 |
| 2 | no no no no | 258 | what are you doing | 173 | what are you doing | 173 | no no no no | 255 |
| 3 | what are you doing | 173 | come on come on | 99 | come on come on | 108 | what are you doing | 173 |
| 4 | i don t think | 131 | what do you mean | 83 | what do you mean | 83 | why do n't you | 114 |
| 5 | why don t you | 114 | i want you to | 82 | i want you to | 82 | do n't know what | 110 |
| 6 | come on come on | 111 | you want me to | 61 | oh oh oh oh | 71 | come on come on | 108 |
| 7 | don t know what | 111 | aargh aargh aargh aargh | 58 | aargh aargh aargh aargh | 68 | i do n't think | 101 |
| 8 | i don t think | 101 | are you talking about | 58 | you want me to | 61 | i do n't want | 95 |
| 9 | i don t want | 95 | oh oh oh oh | 58 | i know i know | 60 | you do n't have | 88 |
| 10 | i m i m | 91 | what do you think | 58 | what do you think | 59 | what do you mean | 83 |

## COLLOCATIONS AND CONCORDANCES

Chronologically, the notion of collocation was first introduced by Palmer (1933) who defined collocation as a succession of two or more words that must be learned as an integral whole and not pieced together from its component parts and, some years later, by Firth (1957b:14) who defined it as "actual words in habitual company". Firth (1951, 1957a) particularly emphasized the habituality which distinguishes collocation and the limited possibility of co-occurrence of words, or, in Sinclairian modern terms, the phraseological tendency of language: "One of the meanings of ass is its habitual collocation with an immediately preceding you silly, and with other phrases of address or of personal reference. [...] There are only limited possibilities of collocation with preceding adjectives, among which the commonest are silly, obstinate, stupid, awful, occasionally egregious" (Firth 1957a:195). Other scholars gave, then, a slightly different definition of collocation: Leech (1974), for example, pointed out the psychological association "a word acquires on account of the meanings of words which tend to occur in its environment" (Leech 1974:20). Sinclair (1991:170), instead, emphasized the textual trait of collocation, i.e. "the occurrence of two or more words within a short space of each other in a text". Both Hoey (1991) and Stubbs (2001) highlighted its statistical aspect, namely the chance of relationship that "a lexical item has with items that appear with greater than random probability in its (textual) context" (Hoey 1991:6-7), or, simply, "frequent co-occurrence" (Stubbs 2001:29). However, despite these different slants (i.e. contextual, psychological, textual, and statistical), what remains at the basis of the notion of collocation is the Firthian intuition that the meaning created by the co-occurrence of two items in a given context is a product of those two co-occurring words in that particular context, or in Hallidayan terms, "of the relationship between the system and its environment" (Halliday 2003:196, cf. also Halliday 1985).
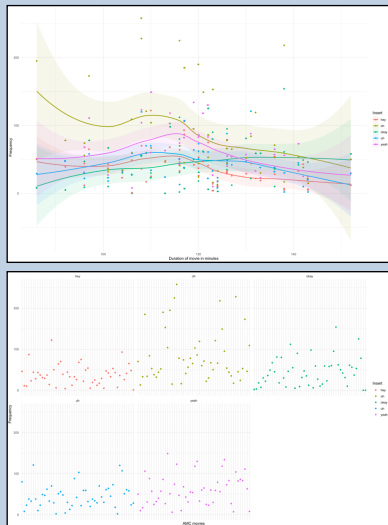
## MULTI-DIMENSIONAL ANALYSIS (MDA)

Technically, via factor analysis a large number of linguistic features characterizing a text are reduced to a small set of derived variables called Factors. Then, through a calculation of the communicative functions most widely shared by the linguistic features in question, each Factor is interpreted functionally as a Dimension of variation which underlines each set of co-occurring linguistic features . More specifically, the following five Biberian dimensions, which are represented by Factors 1-5 respectively, are considered :

D1: the informational (negative) vs. involved (positive) production dimension, which identifies whether a text is marked by high informational density and exact informational content or, on the contrary, by affective, interactional, and generalized content (Biber 1988:107);

D2: the narrative (positive) vs. non-narrative (negative) concerns dimension, which distinguishes narrative discourse from other types of discourse (Biber 1988:109);

D3: the explicit (positive) vs. situation-dependent (negative) reference dimension, which distinguishes between highly explicit, context-independent reference and non-specific, situation-dependent reference (Biber 1988:110);

D4: the overt expression of persuasion (positive) dimension, which marks the degree to which persuasion is marked overtly employed (Biber 1988:111);

D5: the abstract (positive) vs. non-abstract (negative) information dimension, which "seems to mark informational discourse that is abstract, technical, and formal versus other types of discourse" (Biber 1988:113).
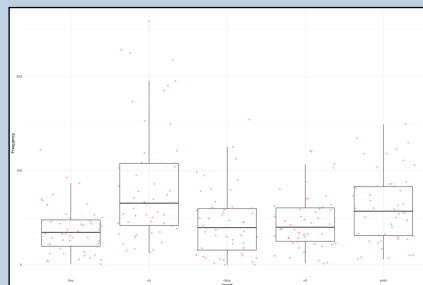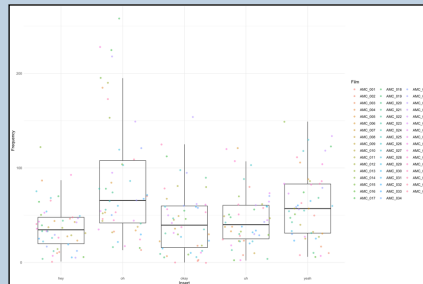
## SCATTERPLOTS

A SCATTERPLOT is "a two-dimensional coordinate system in which the values of the vector are interpreted as coordinates of the y-axis, and the order in which they appear in the vector are the coordinates of the x-axis" (Gries 2009:98). A regression line is a straight line that describes how a response variable y changes as an explanatory variable x changes. We often use a regression line to predict the value of y for a given value of x. Put simply, a scatterplot is a graph which displays a vector according to two dimensions on its x and y axes. A vector is an object with a magnitude and a direction, so that the y axis represents the magnitude, that is the size, while the x axis represents the order in which it appears (i.e. the direction). It is also possible to superimpose a line on the scatterplot which summarizes the relationship between the y and x axes' variables. This line is called "regression line" and it is helpful when we wish to visualize the trend of the vector which will then need to be verified with appropriate tests.
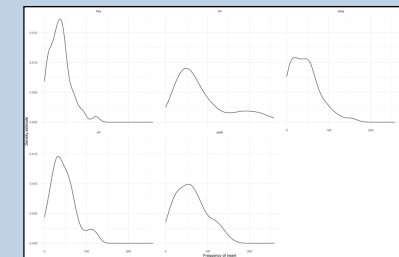




## BOXPLOTS

A BOXPLOT contains various types of valuable information (Gries 2009:119):

• the bold-typed horizontal lines represent the medians of the two vectors;

• the regular horizontal lines that make up the upper and lower boundary of the boxes represent the hinges (approximately the 75%- and the 25% quartiles);

• the whiskers – the dashed vertical lines extending from the box until the upper and lower limit – represent the largest and smallest values that are not more than 1.5 interquartile ranges away from the box;

• each outlier that would be outside of the range of the whiskers would be represented with an individual dot;

• the notches on the left and right sides of the boxes extend across the range ±1.58*IQR/sqrt(n): if the notches of two boxplots overlap, then these will most likely not be significantly different.





## DENSITY PLOTS

A DENSITY PLOT shows the ordered numerical values of a variable x on the horizontal axis, and the probability density of x on the vertical axis (Levshina 2015:51). Put simply, a density plot is a useful graph when one wishes to display the distribution of the data. If the distribution appears to be normal (and this is verified with a statistical test called Shapiro-Wilk test), then parametric tests can be employed to explore the dataset. However, if the density plot shows that the data is skewed, that is, not normally distributed (and this is verified via means of a Shapiro-Wilk test), then the non-parametric version of the tests should be preferred. This graph plots all the numerical values of the data on the x axis, while the y axis corresponds to the probability density. Thus, any peaks in the curve are to be interpreted as follows: the majority of numerical values are distributed along this peak, while the remaining ones are found along its tail(s).