

# Package ‘mda.biber’

October 7, 2025

**Title** Functions for Multi-Dimensional Analysis

**Version** 1.0.1

**Date** 2025-09-22

**Description** Multi-Dimensional Analysis (MDA) is an adaptation of factor analysis developed by Douglas Biber (1992) <[doi:10.1007/BF00136979](https://doi.org/10.1007/BF00136979)>. Its most common use is to describe language as it varies by genre, register, and use. This package contains functions for carrying out the calculations needed to describe and plot MDA results: dimension scores, dimension means, and factor loadings.

**License** MIT + file LICENSE

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 7.3.2

**Imports** dplyr, ggplot2, ggpubr, ggrepel, nFactors, stats, tidyr,  
viridis

**Depends** R (>= 2.10)

**Suggests** rmarkdown, knitr, corrplot, kableExtra, tidyverse, testthat  
(>= 3.0.0)

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** David Brown [aut, cre] (ORCID: <<https://orcid.org/0000-0001-7745-6354>>),  
Alex Reinhart [aut] (ORCID: <<https://orcid.org/0000-0002-6658-514X>>)

**Maintainer** David Brown <dwb2@andrew.cmu.edu>

**Repository** CRAN

**Date/Publication** 2025-10-07 18:00:02 UTC

## Contents

boxplot_mda . . . . .	2
mda_loadings . . . . .	2

micusp_biber . . . . .	4
screeplot_mda . . . . .	6
stickplot_mda . . . . .	7

<b>Index</b>	<b>8</b>
--------------	----------

---

boxplot_mda	<i>Create boxplot for multi-dimensional analysis</i>
-------------	--

---

### Description

Combine scaled vectors of the relevant factor loadings and boxplots of dimension scores.

### Usage

```
boxplot_mda(mda_data, n_factor = 1)
```

### Arguments

mda_data	An mda data.frame produced by the mda_loadings() function.
n_factor	The factor to be plotted.

### Value

A combined plot of scaled vectors and boxplots.

### See Also

[stickplot\\_mda\(\)](#)

---

mda_loadings	<i>Conduct multi-dimensional analysis</i>
--------------	---

---

### Description

Multi-Dimensional Analysis is a statistical procedure developed by Biber and is commonly used in descriptions of language as it varies by genre, register, and task. The procedure is a specific application of factor analysis, which is used as the basis for calculating a 'dimension score' for each text.

### Usage

```
mda_loadings(obs_by_group, n_factors, cor_min = 0.2, threshold = 0.35)
```

## Arguments

obs_by_group	A data frame containing exactly 1 categorical (factor) variable and multiple continuous (numeric) variables. Each row represents one document/observation.
n_factors	The number of factors to be calculated in the factor analysis.
cor_min	The correlation threshold for including variables in the factor analysis. Variables whose (absolute) Pearson correlation with any other variable is greater than this threshold will be included in the factor analysis. Set to 0 to disable thresholding.
threshold	The loading threshold above which variables should be included in factor score calculations. Set to 0 to include all variables.

## Details

MDA is fundamentally factor analysis using the promax rotation, applied to the numeric variables in obs\_by\_group. However, MDA adds two screening steps:

1. Only variables with a nontrivial correlation with any other variable are included; the correlation threshold is configurable with the cor\_min argument.
2. The factor scores are based only on variables whose loadings are greater (in absolute value) than the threshold argument. (Variables are standardized to ensure loadings are comparable.)

These two choices eliminate variables that are uncorrelated with others, and essentially enforce sparsity in each factor, ensuring it is loaded only on a smaller set of variables.

## Value

An mda data frame containing one row per document, containing factor scores for each document. Attributes include the number of factors (n\_factors), the correlation threshold (threshold), the factor loadings (loadings), and the mean factor score for each group (group\_means).

## References

Biber (1988). *Variation across Speech and Writing*. Cambridge University Press.

Biber (1992). "The multi-dimensional approach to linguistic analyses of genre variation: An overview of methodology and findings." *Computers and the Humanities* 26 (5/6), 331-345. doi:10.1007/BF00136979

## See Also

[screplot\\_mda\(\)](#), [stickplot\\_mda\(\)](#), [boxplot\\_mda\(\)](#)

## Examples

```
# Extract the subject area from each document ID and use it as the grouping
# variable
micusp_biber$doc_id <- factor(substr(micusp_biber$doc_id, 1, 3))

m <- mda_loadings(micusp_biber, n_factors = 2)

attr(m, "group_means")
```

heatmap\_mda(m)

---

micusp\_biber

*MICUSP corpus tagged with pseudobibeR features*

---

## Description

The Michigan Corpus of Upper-Level Student Papers (MICUSP) contains 828 student papers. Here each document is tagged with Biber features using the pseudobibeR package. Type-to-token ratio is calculated using the moving average type-to-token ratio (MATTR).

## Usage

micusp\_biber

## Format

A data frame with 828 rows and 68 columns:

**doc\_id** Document ID (from MICUSP)

**f\_01\_past\_tense** Rate of past tense per 1,000 tokens

**f\_02\_perfect\_aspect** Rate of perfect aspect per 1,000 tokens

**f\_03\_present\_tense** Rate of present tense per 1,000 tokens

**f\_04\_place\_adverbials** Rate of place adverbials per 1,000 tokens

**f\_05\_time\_adverbials** Rate of time adverbials per 1,000 tokens

**f\_06\_first\_person\_pronouns** Rate of first person pronouns per 1,000 tokens

**f\_07\_second\_person\_pronouns** Rate of second person pronouns per 1,000 tokens

**f\_08\_third\_person\_pronouns** Rate of third person pronouns per 1,000 tokens

**f\_09\_pronoun\_it** Rate of pronoun 'it' per 1,000 tokens

**f\_10\_demonstrative\_pronoun** Rate of demonstrative pronouns per 1,000 tokens

**f\_11\_indefinite\_pronouns** Rate of indefinite pronouns per 1,000 tokens

**f\_12\_proverb\_do** Rate of proverb 'do' per 1,000 tokens

**f\_13\_wh\_question** Rate of wh-questions per 1,000 tokens

**f\_14\_nominalizations** Rate of nominalizations per 1,000 tokens

**f\_15\_gerunds** Rate of gerunds per 1,000 tokens

**f\_16\_other\_nouns** Rate of other nouns per 1,000 tokens

**f\_17\_agentless\_passives** Rate of agentless passives per 1,000 tokens

**f\_18\_by\_passives** Rate of by-passives per 1,000 tokens

**f\_19\_be\_main\_verb** Rate of 'be' as main verb per 1,000 tokens

**f\_20\_existential\_there** Rate of existential 'there' per 1,000 tokens

**f\_21\_that\_verb\_comp** Rate of that-verb complements per 1,000 tokens  
**f\_22\_that\_adj\_comp** Rate of that-adjective complements per 1,000 tokens  
**f\_23\_wh\_clause** Rate of wh-clauses per 1,000 tokens  
**f\_24\_infinitives** Rate of infinitives per 1,000 tokens  
**f\_25\_present\_participle** Rate of present participles per 1,000 tokens  
**f\_26\_past\_participle** Rate of past participles per 1,000 tokens  
**f\_27\_past\_participle\_whiz** Rate of past participle whiz-deletions per 1,000 tokens  
**f\_28\_present\_participle\_whiz** Rate of present participle whiz-deletions per 1,000 tokens  
**f\_29\_that\_subj** Rate of that-subject clauses per 1,000 tokens  
**f\_30\_that\_obj** Rate of that-object clauses per 1,000 tokens  
**f\_31\_wh\_subj** Rate of wh-subject clauses per 1,000 tokens  
**f\_32\_wh\_obj** Rate of wh-object clauses per 1,000 tokens  
**f\_33\_pied\_piping** Rate of pied-piping per 1,000 tokens  
**f\_34\_sentence\_relatives** Rate of sentence relatives per 1,000 tokens  
**f\_35\_because** Rate of 'because' per 1,000 tokens  
**f\_36\_though** Rate of 'though' per 1,000 tokens  
**f\_37\_if** Rate of 'if' per 1,000 tokens  
**f\_38\_other\_adv\_sub** Rate of other adverbial subordinators per 1,000 tokens  
**f\_39\_prepositions** Rate of prepositions per 1,000 tokens  
**f\_40\_adj\_attr** Rate of attributive adjectives per 1,000 tokens  
**f\_41\_adj\_pred** Rate of predicative adjectives per 1,000 tokens  
**f\_42\_adverbs** Rate of adverbs per 1,000 tokens  
**f\_43\_type\_token** Type-token ratio (MATTR)  
**f\_44\_mean\_word\_length** Mean word length  
**f\_45\_conjuncts** Rate of conjuncts per 1,000 tokens  
**f\_46\_downtoners** Rate of downtoners per 1,000 tokens  
**f\_47\_hedges** Rate of hedges per 1,000 tokens  
**f\_48\_amplifiers** Rate of amplifiers per 1,000 tokens  
**f\_49\_emphatics** Rate of emphatics per 1,000 tokens  
**f\_50\_discourse\_particles** Rate of discourse particles per 1,000 tokens  
**f\_51\_demonstratives** Rate of demonstratives per 1,000 tokens  
**f\_52\_modal\_possibility** Rate of possibility modals per 1,000 tokens  
**f\_53\_modal\_necessity** Rate of necessity modals per 1,000 tokens  
**f\_54\_modal\_predictive** Rate of predictive modals per 1,000 tokens  
**f\_55\_verb\_public** Rate of public verbs per 1,000 tokens  
**f\_56\_verb\_private** Rate of private verbs per 1,000 tokens  
**f\_57\_verb\_suasive** Rate of suasive verbs per 1,000 tokens

**f\_58\_verb\_seem** Rate of 'seem' verbs per 1,000 tokens  
**f\_59\_contractions** Rate of contractions per 1,000 tokens  
**f\_60\_that\_deletion** Rate of that-deletions per 1,000 tokens  
**f\_61\_stranded\_preposition** Rate of stranded prepositions per 1,000 tokens  
**f\_62\_split\_infinitive** Rate of split infinitives per 1,000 tokens  
**f\_63\_split\_auxiliary** Rate of split auxiliaries per 1,000 tokens  
**f\_64\_phrasal\_coordination** Rate of phrasal coordination per 1,000 tokens  
**f\_65\_clausal\_coordination** Rate of clausal coordination per 1,000 tokens  
**f\_66\_neg\_synthetic** Rate of synthetic negation per 1,000 tokens  
**f\_67\_neg\_analytic** Rate of analytic negation per 1,000 tokens

### Source

Michigan Corpus of Upper-Level Student Papers, <https://elicorpora.info/main>, tagged with the pseudobiber package.

---

screepLOT\_mda

*Scree plot for multi-dimensional analysis*

---

### Description

The scree plot shows each factor along the X axis, and the proportion of common variance explained by that factor on the Y axis. The proportion of common variance explained is given by the factor eigenvalue.

### Usage

```
screepLOT_mda(obs_by_group, cor_min = 0.2)
```

### Arguments

**obs\_by\_group** A data frame containing 1 categorical (factor) variable and continuous (numeric) variables.  
**cor\_min** The correlation threshold for including variables in the factor analysis.

### Details

A wrapper for the `nFactors::nScree()` and `nFactors::plotnScree()` functions.

### Value

Nothing returned

### See Also

[mda\\_loadings\(\)](#)

---

stickplot_mda	<i>Plots of MDA factor group means and loadings</i>
---------------	---

---

**Description**

Stick plots show each group's mean loading along a factor, plotted along a positive/negative cline. Heatmaps show each variable's loading on a factor. `stickplot_mda()` produces just a stick plot, while `heatmap_mda()` places a heatmap alongside the stick plot.

**Usage**

```
stickplot_mda(mda_data, n_factor = 1)
```

```
heatmap_mda(mda_data, n_factor = 1)
```

**Arguments**

<code>mda_data</code>	An mda data frame produced by the <code>mda_loadings()</code> function.
<code>n_factor</code>	Index of the factor to be plotted.

**Value**

ggplot object

**See Also**

[boxplot\\_mda\(\)](#)

# Index

## \* datasets

micusp\_biber, 4

boxplot\_mda, 2

boxplot\_mda(), 3, 7

heatmap\_mda (stickplot\_mda), 7

mda\_loadings, 2

mda\_loadings(), 6

micusp\_biber, 4

screeplot\_mda, 6

screeplot\_mda(), 3

stickplot\_mda, 7

stickplot\_mda(), 2, 3