

# Package ‘docextractr’

January 9, 2019

**Title** Extract Data Tables and Comments from 'Microsoft' 'Word' Documents

**Version** 0.6.1

**Maintainer** Bob Rudis <bob@rud.is>

**Description** 'Microsoft Word' 'docx' files provide an 'XML' structure that is fairly straightforward to navigate, especially when it applies to 'Word' tables and comments. Tools are provided to determine table count/structure, comment count and also to extract/clean tables and comments from 'Microsoft Word' 'docx' documents. There is also nascent support for '.doc' files.

**SystemRequirements** LibreOffice (<<https://www.libreoffice.org/>>) required to extract data from .doc files.

**URL** <http://gitlab.com/hrbrmstr/docextractr>

**BugReports** <https://gitlab.com/hrbrmstr/docextractr/issues>

**Encoding** UTF-8

**Depends** R (>= 3.2.0)

**License** MIT + file LICENSE

**LazyData** true

**Suggests** testthat, covr

**Imports** tools, xml2, purrr, dplyr, utils, httr, magrittr

**RoxygenNote** 6.1.1

**NeedsCompilation** no

**Author** Bob Rudis [aut, cre] (<<https://orcid.org/0000-0001-5670-2640>>),  
Mark Dulhunty [ctb],  
Karlo Guidoni-Martins [ctb],  
Chris Muir [aut, ctb]

**Repository** CRAN

**Date/Publication** 2019-01-09 17:00:10 UTC

## R topics documented:

assign_colnames . . . . .	2
docextractr . . . . .	3
docx_cmnt_count . . . . .	3
docx_describe_cmnts . . . . .	4
docx_describe_tbls . . . . .	4
docx_extract_all . . . . .	5
docx_extract_all_cmnts . . . . .	6
docx_extract_all_tbls . . . . .	6
docx_extract_tbl . . . . .	7
docx_tbl_count . . . . .	8
mcga . . . . .	8
print.docx . . . . .	9
read_docx . . . . .	9
set_libreoffice_path . . . . .	10

<b>Index</b>	<b>12</b>
--------------	-----------

---

assign_colnames	<i>Make a specific row the column names for the specified data.frame</i>
-----------------	--

---

### Description

Many tables in Word documents are in twisted formats where there may be labels or other oddities mixed in that make it difficult to work with the underlying data. This function makes it easy to identify a particular row in a scraped data.frame as the one containing column names and have it become the column names, removing it and (optionally) all of the rows before it (since that's usually what needs to be done).

### Usage

```
assign_colnames(dat, row, remove = TRUE, remove_previous = remove)
```

### Arguments

dat	can be any data.frame but is intended for use with ones returned by this package
row	numeric value indicating the row number that is to become the column names
remove	remove row specified by row after making it the column names? (Default: TRUE)
remove_previous	remove any rows preceding row? (Default: TRUE but will be assigned whatever is given for remove).

### Value

data.frame

**See Also**

[docx\\_extract\\_all](#), [docx\\_extract\\_tbl](#)

**Examples**

```
# a "real" Word doc
real_world <- read_docx(system.file("examples/realworld.docx", package="docxtractr"))
docx_tbl_count(real_world)

# get all the tables
tbls <- docx_extract_all_tbls(real_world)

# make table 1 better
assign_colnames(tbls[[1]], 2)

# make table 5 better
assign_colnames(tbls[[5]], 2)
```

---

docxtractr	<i>Extract Data Tables and Comments from 'Microsoft' 'Word' Documents</i>
------------	---

---

**Description**

Microsoft Word ‘docx’ files provide an XML structure that is fairly straightforward to navigate, especially when it applies to Word tables. The ‘docxtractr’ package provides tools to determine table count + table structure and extract tables from Microsoft Word docx documents. It also provides tools to determine comment count and extract comments from Word ‘docx’ documents.

**Author(s)**

Bob Rudis (bob@rud.is)

---

docx_cmnt_count	<i>Get number of comments in a Word document</i>
-----------------	--

---

**Description**

Get number of comments in a Word document

**Usage**

```
docx_cmnt_count(docx)
```

**Arguments**

docx	docx object read with read_docx
------	---------------------------------

**Value**

numeric

**Examples**

```
cmnts <- read_docx(system.file("examples/comments.docx", package="docxtractr"))
docx_cmnt_count(cmnts)
```

---

docx\_describe\_cmnts *Returns information about the comments in the Word document*

---

**Description**

Returns information about the comments in the Word document

**Usage**

```
docx_describe_cmnts(docx)
```

**Arguments**

docx                    docx object read with read\_docx

**Examples**

```
cmnts <- read_docx(system.file("examples/comments.docx", package="docxtractr"))
docx_cmnt_count(cmnts)
docx_describe_cmnts(cmnts)
```

---

docx\_describe\_tbls *Returns a description of all the tables in the Word document*

---

**Description**

This function will attempt to discern the structure of each of the tables in docx and print this information

**Usage**

```
docx_describe_tbls(docx)
```

**Arguments**

docx                    docx object read with read\_docx

## Examples

```
complx <- read_docx(system.file("examples/complex.docx", package="docxtractr"))
docx_tbl_count(complx)
docx_describe_tbls(complx)
```

---

docx_extract_all	<i>Extract all tables from a Word document</i>
------------------	--

---

## Description

Extract all tables from a Word document

## Usage

```
docx_extract_all(docx, guess_header = TRUE, preserve = FALSE,
  trim = TRUE)
```

## Arguments

docx	docx object read with read_docx
guess_header	should the function make a guess as to the existence of a header in a table? (Default: TRUE)
preserve	preserve line breaks within a cell? Default: 'FALSE'. NOTE: This overrides 'trim'.
trim	trim leading/trailing whitespace (if any) in cells? (default: TRUE)

## Value

list of data.frames or an empty list if no tables exist in docx

## See Also

[assign\\_colnames](#), [docx\\_extract\\_tbl](#)

## Examples

```
# a "real" Word doc

real_world <- read_docx(system.file("examples/realworld.docx", package="docxtractr"))
docx_tbl_count(real_world)

# get all the tables
tbls <- docx_extract_all_tbls(real_world)
```

---

docx\_extract\_all\_cmnts

*Extract all comments from a Word document*

---

### Description

Extract all comments from a Word document

### Usage

```
docx_extract_all_cmnts(docx, include_text = FALSE)
```

### Arguments

docx                docx object read with read\_docx  
include\_text        if TRUE then the text associated with the comment will also be included

### Value

data\_frame of comment id, author & text

### Examples

```
cmnts <- read_docx(system.file("examples/comments.docx", package="docxtractr"))
docx_cmnt_count(cmnts)
docx_describe_cmnts(cmnts)
docx_extract_all_cmnts(cmnts)
```

---

docx\_extract\_all\_tbls *Extract all tables from a Word document*

---

### Description

Extract all tables from a Word document

### Usage

```
docx_extract_all_tbls(docx, guess_header = TRUE, preserve = FALSE,
  trim = TRUE)
```

### Arguments

docx                docx object read with read\_docx  
guess\_header        should the function make a guess as to the existence of a header in a table?  
                      (Default: TRUE)  
preserve            preserve line breaks within a cell? Default: 'FALSE'. NOTE: This overrides  
                      'trim'.  
trim                 trim leading/trailing whitespace (if any) in cells? (default: TRUE)

**Value**

list of data.frames or an empty list if no tables exist in docx

**See Also**

[assign\\_colnames](#), [docx\\_extract\\_tbl](#)

**Examples**

```
# a "real" Word doc

real_world <- read_docx(system.file("examples/realworld.docx", package="docxtractr"))
docx_tbl_count(real_world)

# get all the tables
tbls <- docx_extract_all_tbls(real_world)
```

---

docx_extract_tbl	<i>Extract a table from a Word document</i>
------------------	---

---

**Description**

Given a document read with `read_docx` and a table to extract (optionally indicating whether there was a header or not and if cell whitespace trimming is desired) extract the contents of the table to a `data.frame`.

**Usage**

```
docx_extract_tbl(docx, tbl_number = 1, header = TRUE,
  preserve = FALSE, trim = TRUE)
```

**Arguments**

<code>docx</code>	docx object read with <code>read_docx</code>
<code>tbl_number</code>	which table to extract (defaults to 1)
<code>header</code>	assume first row of table is a header row? (default: TRUE)
<code>preserve</code>	preserve line breaks within a cell? Default: FALSE. NOTE: This overrides <code>trim</code> .
<code>trim</code>	trim leading/trailing whitespace (if any) in cells? (default: TRUE)

**Value**

`data.frame`

**See Also**

[docx\\_extract\\_all](#), [docx\\_extract\\_tbl](#), [assign\\_colnames](#)

**Examples**

```
doc3 <- read_docx(system.file("examples/data3.docx", package="docxtractr"))
docx_extract_tbl(doc3, 3)

intracell_whitespace <- read_docx(system.file("examples/preserve.docx", package="docxtractr"))
docx_extract_tbl(intracell_whitespace, 2, preserve=FALSE)
docx_extract_tbl(intracell_whitespace, 2, preserve=TRUE)
```

---

docx_tbl_count	<i>Get number of tables in a Word document</i>
----------------	--

---

**Description**

Get number of tables in a Word document

**Usage**

```
docx_tbl_count(docx)
```

**Arguments**

docx	docx object read with read_docx
------	---------------------------------

**Value**

numeric

**Examples**

```
complx <- read_docx(system.file("examples/complex.docx", package="docxtractr"))
docx_tbl_count(complx)
```

---

mcga	<i>Make Column Names Great Again</i>
------	--------------------------------------

---

**Description**

Remove punctuation and spaces and turn them to underscores plus convert to lower case.

**Usage**

```
mcga(tbl)
```

**Arguments**

tbl	a data.frame-like object
-----	--------------------------



**Value**

whatever class `x` was but with truly great, really great column names. They're amazing. Trust me. They'll be incredible column names once we're done.

**Examples**

```
real_world <- read_docx(system.file("examples/realworld.docx", package="docxtractr"))
tbls <- docx_extract_all_tbls(real_world)
mcga(assign_colnames(tbls[[1]], 2))
```

---

print.docx

*Display information about the document*


---

**Description**

Display information about the document

**Usage**

```
## S3 method for class 'docx'
print(x, ...)
```

**Arguments**

<code>x</code>	docx object
<code>...</code>	ignored

---

read\_docx

*Read in a Word document for table extraction*


---

**Description**

Local file path or URL pointing to a .docx file. Can also take .doc file as input if LibreOffice is installed (see <https://www.libreoffice.org/> for more info and to download).

**Usage**

```
read_docx(path, track_changes = NULL)
```

**Arguments**

path	path to the Word document
track_changes	if not NULL (the default) then must be one of "accept" or "reject" which will, respectively, accept all or reject all changes. NOTE: this functionality relies on the pandoc utility being available on the system PATH. Both system PATH and the RSTUDIO_PANDOC (RStudio ships with a copy of pandoc) environment variables will be checked. If no pandoc binary is found then a warning will be issued and the document will be read without integrating or ignoring any tracked changes. The original Word document <i>will not be modified</i> and this feature <i>only works</i> with docx files.

**Examples**

```
doc <- read_docx(system.file("examples/data.docx", package="docxtractr"))
class(doc)

doc <- read_docx(
  system.file("examples/trackchanges.docx", package="docxtractr"),
  track_changes = "accept"
)

## Not run:
# from a URL

## End(Not run)
```

---

set\_libreoffice\_path *Point to Local soffice.exe File*

---

**Description**

Function to set an option that points to the local LibreOffice file soffice.exe.

**Usage**

```
set_libreoffice_path(path)
```

**Arguments**

path	path to the LibreOffice soffice file
------	--------------------------------------

**Details**

For a list of possible file path locations for soffice.exe, see <https://github.com/hrbrmstr/docxtractr/issues/5#issuecomment-233181976>

**Value**

Returns nothing, function sets the option variable `path_to_libreoffice`.

**Examples**

```
## Not run:  
set_libreoffice_path("local/path/to/soffice.exe")  
  
## End(Not run)
```

# Index

`assign_colnames`, [2](#), [5](#), [7](#)

`docx_cmnt_count`, [3](#)

`docx_describe_cmnts`, [4](#)

`docx_describe_tbls`, [4](#)

`docx_extract_all`, [3](#), [5](#), [7](#)

`docx_extract_all_cmnts`, [6](#)

`docx_extract_all_tbls`, [6](#)

`docx_extract_tbl`, [3](#), [5](#), [7](#), [7](#)

`docx_tbl_count`, [8](#)

`docxtractr`, [3](#)

`docxtractr-package (docxtractr)`, [3](#)

`mcga`, [8](#)

`print.docx`, [9](#)

`read_docx`, [9](#)

`set_libreoffice_path`, [10](#)