

# Package: regextable (via r-universe)

May 20, 2026

**Title** Pattern-Based Text Extraction and Standardization with Lookup Tables

**Version** 0.1.2

**Description** Extracts information from text using lookup tables of regular expressions. Each text entry is compared against patterns in a lookup table, returning the extracted substrings alongside optional metadata. Users can choose to extract all matching patterns (generating multiple rows per text entry) or limit extraction to the first match. This approach enables comprehensive, standardized pattern coverage when processing large or complex text datasets.

**LazyData** true

**License** MIT + file LICENSE

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.3

**Imports** chk, dplyr, stringi, stringr, pbapply, stats

**Suggests** kableExtra, knitr, rmarkdown, quarto, spelling, testthat (>= 3.0.0), spacyr, tibble, googledrive, googlesheets4

**VignetteBuilder** quarto

**URL** <https://github.com/judgelord/regextable>,  
<https://judgelord.github.io/regextable/>

**BugReports** <https://github.com/judgelord/regextable/issues>

**Config/testthat/edition** 3

**Depends** R (>= 4.1)

**Language** en-US

**Config/pak/sysreqs** libicu-dev

**Repository** <https://judgelord.r-universe.dev>

**Date/Publication** 2026-04-20 19:03:11 UTC

**RemoteUrl** https://github.com/judgelord/regextable

**RemoteRef** HEAD

**RemoteSha** 035ef79284faf43743f3c9fa87c8bb763d299eb2

## Contents

clean_text . . . . .	2
cr2007_03_01 . . . . .	3
extract . . . . .	3
members . . . . .	6

<b>Index</b>	<b>7</b>
--------------	----------

---

clean_text	<i>Clean Text</i>
------------	-------------------

---

## Description

Cleans a character vector by converting to lowercase, removing or replacing specific punctuation, normalizing commas, and squishing excess whitespace.

## Usage

```
clean_text(text)
```

## Arguments

text            Character vector to clean.

## Value

A cleaned character vector.

## Examples

```
clean_text(c("Hello World!?", "This--is\\tR.\\nTesting: 1, 2, , 3;"))
```

---

cr2007_03_01	<i>cr2007_03_01 dataset</i>
--------------	-----------------------------

---

**Description**

Sample text dataset used for demonstration of `regextable`.

**Format**

A tibble with 5 columns:

**date** Date of the record (YYYY-MM-DD)

**speaker** Speaker name in the text

**header** Header or title of the speech

**url** Original URL of the source text

**url\_txt** Full text content from the source

**Source**

Generated for the `regextable` package.

---

extract	<i>Extract Regex Pattern Matches from Text Data</i>
---------	---

---

**Description**

Matches text against a table of regular expressions and returns extracted matches with optional metadata.

**Usage**

```
extract(  
  data,  
  regex_table,  
  col_name = "text",  
  pattern_col = "pattern",  
  typo_table = NULL,  
  typo_from_col = "typo",  
  typo_to_col = "correction",  
  date_col = NULL,  
  date_start = NULL,  
  date_end = NULL,  
  data_return_cols = NULL,  
  regex_return_cols = NULL,
```

```

remove_acronyms = FALSE,
do_clean_text = TRUE,
unique_match = FALSE,
use_ner = FALSE,
ner_timing = "after",
ner_entity_types = c("ORG"),
verbose = TRUE,
cl = NULL
)

```

## Arguments

<code>data</code>	A data frame or character vector containing the text to search. If a character vector is provided, it is internally converted to a data frame and <code>col_name</code> is ignored.
<code>regex_table</code>	A data frame containing regular expression patterns and optional metadata columns.
<code>col_name</code>	Character string specifying the column in <code>data</code> that contains text to search. Default is "text".
<code>pattern_col</code>	Character string specifying the column in <code>regex_table</code> that contains regex patterns. Default is "pattern".
<code>typo_table</code>	Optional data frame with text replacements for corrections. Replacements are applied sequentially to the text using regex (with word boundaries) before pattern matching.
<code>typo_from_col</code>	Optional column in <code>typo_table</code> with text to replace. Default is "typo".
<code>typo_to_col</code>	Optional column in <code>typo_table</code> with replacement text. Default is "correction".
<code>date_col</code>	Optional column in <code>data</code> for date filtering. If provided, rows are filtered by <code>date_start</code> and <code>date_end</code> before pattern matching.
<code>date_start</code>	Optional start date (Date object or string like "YYYY-MM-DD") for filtering data when <code>date_col</code> is specified.
<code>date_end</code>	Optional end date (Date object or string like "YYYY-MM-DD") for filtering data when <code>date_col</code> is specified.
<code>data_return_cols</code>	Optional vector of column names to include from <code>data</code> . Default is NULL (only <code>row_id</code> is returned).
<code>regex_return_cols</code>	Optional vector of column names to include from <code>regex_table</code> . Default is NULL (no metadata columns added).
<code>remove_acronyms</code>	Logical; if TRUE, removes patterns consisting only of uppercase letters (2 or more characters) from <code>regex_table</code> .
<code>do_clean_text</code>	Logical; if TRUE, applies basic text cleaning to the input before matching.
<code>unique_match</code>	Logical; if TRUE, stops searching after the first match to find at most one match per row (evaluated in the order patterns appear in <code>regex_table</code> ). If FALSE, returns all matches for all patterns.

<code>use_ner</code>	Logical; if TRUE, uses the 'spacyr' package to validate that matches are actual Named Entities (e.g., organizations). Note: spacyr must be initialized (e.g., via <code>spacyr::spacy_initialize()</code> ) before calling this function.
<code>ner_timing</code>	Character string; either "after" or "before". If "after" (default), regex matches are found first, then validated with NER. If "before", NER extracts entities first, and regex searches only within those entities.
<code>ner_entity_types</code>	Character vector; the types of Named Entities to keep if <code>use_ner</code> is TRUE. Default is "ORG".
<code>verbose</code>	Logical; if TRUE, displays progress messages.
<code>cl</code>	A cluster object created by <code>parallel::makeCluster()</code> , or an integer to indicate number of child processes (integer values are ignored on Windows). Passed to <code>pbapply::pblapply()</code> .

### Details

Pattern matching is performed using R's regular expression engine and is case-insensitive by default. For each input row, the function checks patterns in `regex_table` and returns matches based on the `unique_match` parameter.

### Value

A tibble with the following columns:

- `row_id`: Integer identifier corresponding to rows in the input data.
- Additional columns from data if `data_return_cols` is specified.
- Additional columns from `regex_table` if `regex_return_cols` is specified.
- `pattern`: The matched regular expression pattern(s).
- `match`: The extracted text from the data (original casing preserved).

### Examples

```
# Create sample data
data <- data.frame(
  id = 1:3,
  text = c("I love apples", "Bananas are great", "Oranges and apples"),
  stringsAsFactors = FALSE
)

# Create regex patterns
patterns <- data.frame(
  pattern = c("apples", "bananas", "oranges"),
  category = c("fruit", "fruit", "fruit")
)

# Extract all matches
extract(data, patterns)

# Extract one match per row
extract(data, patterns, unique_match = TRUE)
```

---

members

*members dataset*

---

### **Description**

Lookup table of member names and metadata for regex matching.

### **Format**

A tibble with 9 columns:

**congress** Congress number (numeric)

**chamber** Chamber (House/President/Senate)

**bioname** Full bio name of the member

**pattern** Regex pattern to match this member's name

**icpsr** Numeric ICPSR identifier

**state\_abbrev** Two-letter state abbreviation

**district\_code** District number (0 for President)

**first\_name** First name of the member

**last\_name** Last name of the member

### **Source**

Generated for the `regextable` package.

# Index

`clean_text`, 2  
`cr2007_03_01`, 3

`extract`, 3

`members`, 6

`pbapply::pblapply()`, 5