

# Package ‘proustr’

October 4, 2017

**Title** Tools for Natural Language Processing in French

**Version** 0.2.1

**Date** 2017-09-28

## Description

Tools for Natural Language Processing in French and texts from Marcel Proust's collection ``A La Recherche Du Temps Perdu". The novels contained in this collection are ``Du cote de chez Swann ", ``A l'ombre des jeunes filles en fleurs", ``Le Cote de Guermantes", ``Sodome et Gomorrhe I et II", ``La Prisonniere", ``Albertine disparue", and ``Le Temps retrouve".

**URL** <https://github.com/ColinFay/proustr>

**BugReports** <https://github.com/ColinFay/proustr/issues>

**Depends** R (>= 2.10)

**License** MIT + file LICENSE

**Imports** dplyr, magrittr, stringr, rlang, purrr, tidyr, tokenizers,  
SnowballC, assertthat

**LazyData** true

**RoxygenNote** 6.0.1.9000

**Encoding** UTF-8

**Suggests** testthat, knitr, rmarkdown, covr

**VignetteBuilder** knitr

**ByteCompile** true

**NeedsCompilation** no

**Author** Colin FAY [aut, cre]

**Maintainer** Colin FAY <contact@colinfay.me>

**Repository** CRAN

**Date/Publication** 2017-10-04 20:15:20 UTC

**R topics documented:**

albertinedisparue	2
alombredesjeunesfillesenfleurs	3
ducotedechezswann	3
laprisonniere	4
lecotedeguermantes	4
letempretrouve	4
proust_books	5
proust_char	5
proust_characters	6
proust_random	6
proust_sentiments	7
proust_stopwords	7
pr_detect_days	8
pr_detect_months	9
pr_detect_pro	9
pr_normalize_punc	10
pr_stem_sentences	11
pr_stem_words	11
sentiments_polarity	12
sentiments_score	12
sodomeetgomorrhe	13
<b>Index</b>	<b>14</b>

---

albertinedisparue	<i>Marcel Proust's novel "Albertine disparue"</i>
-------------------	---

---

**Description**

A dataset containing Marcel Proust's "Albertine disparue". This text has been downloaded from WikiSource.

**Usage**

```
albertinedisparue
```

**Format**

A tibble with text, book, volume, and year

**Source**

[https://fr.wikisource.org/wiki/Albertine\\_disparue](https://fr.wikisource.org/wiki/Albertine_disparue)

---

alombredesjeunesfillesenfleurs

*Marcel Proust's novel "À l'ombre des jeunes filles en fleurs"*

---

### **Description**

A dataset containing Marcel Proust's "À l'ombre des jeunes filles en fleurs". This text has been downloaded from WikiSource.

### **Usage**

alombredesjeunesfillesenfleurs

### **Format**

A tibble with text, book, volume, and year

---

ducotedechezswann

*Marcel Proust's novel "Du côté de chez Swann"*

---

### **Description**

A dataset containing Marcel Proust's "Du côté de chez Swann". This text has been downloaded from WikiSource.

### **Usage**

ducotedechezswann

### **Format**

A tibble with text, book, volume, and year

---

laprisonniere      *Marcel Proust's novel "La Prisonnière"*

---

**Description**

A dataset containing Marcel Proust's "La prisonnière". This text has been downloaded from WikiSource.

**Usage**

laprisonniere

**Format**

A tibble with text, book, volume, and year

---

lecotedeguermantes      *Marcel Proust's novel "Le côté de Guermantes"*

---

**Description**

A dataset containing Marcel Proust's "À l'ombre des jeunes filles en fleurs". This text has been downloaded from WikiSource.

**Usage**

lecotedeguermantes

**Format**

A tibble with text, book, volume, and year

---

letempretrouve      *Marcel Proust's novel "Le temps retrouvé"*

---

**Description**

A dataset containing Marcel Proust's "Le temps retrouvé". This text has been downloaded from WikiSource.

**Usage**

letempretrouve

**Format**

A tibble with text, book, volume, and year.

---

proust\_books

*Tidy data frame of Marcel Proust's 7 novels from La Recherche*

---

### Description

Returns a tidy tibble of Marcel Proust's 7 novels from *À la recherche du temps perdu*. The tibble contains four columns: text, book, volume and year.

### Usage

```
proust_books()
```

### Value

A tibble with four columns: text, book, volume and year.

### Examples

```
#Create the tibble  
proust <- proust_books()
```

---

proust\_char

*Characters from "À la recherche du temps perdu"*

---

### Description

A dataset containing Marcel Proust's characters from "*À la recherche du temps perdu*" and their frequency in each book. This dataset has been downloaded from proust-personnages.

### Usage

```
proust_char
```

### Format

A tibble with their name

### Source

[http://proust-personnages.fr/?page\\_id=10254](http://proust-personnages.fr/?page_id=10254)

---

proust\_characters      *Characters from Proust Books*

---

**Description**

Returns a tidy data frame of Marcel Proust's characters.

**Usage**

```
proust_characters()
```

**Value**

A tibble

**Source**

<http://proust-personnages.fr/>

**Examples**

```
#Creates the tibble  
proust <- proust_characters()
```

---

proust\_random      *Create a Random Proust extract*

---

**Description**

Create your own flavor of Proust with this random extractor.

**Usage**

```
proust_random(count = 1, collapse = TRUE)
```

**Arguments**

count                  the number of line you want to randomly extract and paste.  
collapse                if FALSE, the output will be a tibble. Default is TRUE, a character vector.

**Value**

a character vector

**Examples**

```
proust_random(4)
```

---

*proust\_sentiments*      *Sentiment Lexicon*

---

**Description**

A sentiment lexicon with either polarity or score.

**Usage**

```
proust_sentiments(type = c("polarity", "score"))
```

**Arguments**

`type`                      polarity (positive or negative) or score on six sentiments (joy, fear, sadness, anger, surprise, disgust)

**Value**

a tibble

**Source**

Amine Abdaoui, Jérôme Azé, Sandra Bringay et Pascal Poncelet. FEEL: French Expanded Emotion Lexicon. Language Resources and Evaluation, LRE 2016, pp 1-23.

**Examples**

```
proust_sentiments(type = "score")  
proust_sentiments(type = "polarity")
```

---

*proust\_stopwords*      *Stop Words*

---

**Description**

Stop words concatenated from various web sources.

**Usage**

```
proust_stopwords()
```

**Value**

a tibble with stopwords

**Source**

<https://raw.githubusercontent.com/stopwords-iso/stopwords-fr/master/stopwords-fr.txt>

<http://www.ranks.nl/stopwords/french>

<http://www.naunaute.com/liste-stop-words-francais-393>

[https://sites.google.com/site/kevinbouge/stopwords-lists/stopwords\\_fr.txt?attredirects=0&d=1](https://sites.google.com/site/kevinbouge/stopwords-lists/stopwords_fr.txt?attredirects=0&d=1)

**Examples**

```
proust_stopwords()
```

---

pr_detect_days	<i>Detect french days</i>
----------------	---------------------------

---

**Description**

Detect the name of the days (in French)

**Usage**

```
pr_detect_days(df, col)
```

**Arguments**

df	a dataframe
col	the column containing the text

**Value**

a tibble with the number of days detected by the algo

**Examples**

```
a <- data.frame(jours = c("C'est lundi 1er mars et mardi 2",  
"Et mercredi 3", "Il est revenu jeudi."))  
pr_detect_days(a, jours)
```



---

pr_detect_months	<i>Detect french months</i>
------------------	-----------------------------

---

**Description**

Detect the name of the months (in French)

**Usage**

```
pr_detect_months(df, col)
```

**Arguments**

df	a dataframe
col	the column containing the text

**Value**

a tibble with the number of days detected by the algo

**Examples**

```
a <- data.frame(month = c("C'est lundi 1er mars et mardi 2",  
"Et mercredi 3", "Il est revenu en juin."))  
pr_detect_months(a, month)
```

---

pr_detect_pro	<i>Detect French pronouns</i>
---------------	-------------------------------

---

**Description**

Detect the pronouns from a text (in French)

**Usage**

```
pr_detect_pro(df, col, verbose = FALSE)
```

**Arguments**

df	a dataframe
col	the column containing the text
verbose	wether or not to return the list of pronouns. Defaults is FALSE

**Details**

The shortcuts in the pronoun col stand for:

pps: first person singular (première personne du singulier)

dps: second person singular (deuxième personne du singulier)

tps: third person singular (troisième personne du singulier)

ppp: first person plural (première personne du pluriel)

dpp: second person singular (deuxième personne du pluriel)

tppl: third person singular (troisième personne du pluriel)

**Value**

a tibble with the detected pronouns

**Examples**

```
library(proustr)
a <- proust_books()[1,]
pr_detect_pro(a, text, verbose = TRUE)
pr_detect_pro(a, text)
```

---

pr_normalize_punc	<i>Normalize punctuation</i>
-------------------	------------------------------

---

**Description**

Normalize a text written with usual french punctuation

**Usage**

```
pr_normalize_punc(df, col)
```

**Arguments**

df	a dataframe
col	the column to normalize

**Value**

a tibble with normalized text

**Examples**

```
a <- proustr::albertinedisparue[1:20,]
pr_normalize_punc(albertinedisparue, text)
```

---

pr\_stem\_sentences      *Stem a dataframe containing a column with sentences*

---

**Description**

Implementation of the SnowballC stemmer. Note that punctuation and capital letters are removed when processing.

**Usage**

```
pr_stem_sentences(df, col, language = "french")
```

**Arguments**

df	the data.frame containing the text
col	the column with the text
language	the language of the text. Default is french. See SnowballC::getStemLanguages() function for a list of supported languages.

**Value**

a tibble

**Examples**

```
a <- proustr::laprisonniere[1:10,]  
pr_stem_sentences(a, text)
```

---

pr\_stem\_words      *Stem a dataframe containing a column with words*

---

**Description**

Implementation of the SnowballC stemmer. Note that punctuation and capitals letters are also removed.

**Usage**

```
pr_stem_words(df, col, language = "french")
```

**Arguments**

df	the data.frame containing the sentences
col	the column with the sentences
language	the language of the words Default is french. See SnowballC::getStemLanguages() function for a list of supported languages.

**Value**

a tibble

**Examples**

```
a <- data.frame(words = c("matin", "heure", "fatigué", "sonné", "lois", "tests", "fusionner"))
pr_stem_words(a, words)
```

---

sentiments\_polarity     *Sentiment lexicon with polarity*

---

**Description**

Dataset from: Amine Abdaoui, Jérôme Azé, Sandra Bringay et Pascal Poncelet. FEEL: French Expanded Emotion Lexicon. Language Resources and Evaluation, LRE 2016, pp 1-23.

**Usage**

sentiments\_polarity

**Format**

A tibble

**Source**

<http://www.lirmm.fr/~abdaoui/FEEL.html>

---

sentiments\_score     *Sentiment lexicon with score*

---

**Description**

Dataset from: Amine Abdaoui, Jérôme Azé, Sandra Bringay et Pascal Poncelet. FEEL: French Expanded Emotion Lexicon. Language Resources and Evaluation, LRE 2016, pp 1-23.

**Usage**

sentiments\_score

**Format**

A tibble

**Source**

<http://www.lirmm.fr/~abdaoui/FEEL.html>

---

sodomeetgomorrhe	<i>Marcel Proust's novel "Sodome et Gomorrhe"</i>
------------------	---

---

**Description**

A dataset containing Marcel Proust's "Sodom et Gomorrhe". This text has been downloaded from WikiSource.

**Usage**

sodomeetgomorrhe

**Format**

A tibble with text, book, volume, and year

**Source**

[https://fr.wikisource.org/wiki/Sodome\\_et\\_Gomorrhe](https://fr.wikisource.org/wiki/Sodome_et_Gomorrhe)

# Index

## \*Topic **datasets**

- albertinedisparue, 2
- alombredesjeunesfillesenfleurs, 3
- ducotedechezswann, 3
- laprisonniere, 4
- lecotedeguermantes, 4
- letempretrouve, 4
- proust\_char, 5
- sentiments\_polarity, 12
- sentiments\_score, 12
- sodomeetgomorrhe, 13

  

- albertinedisparue, 2
- alombredesjeunesfillesenfleurs, 3

  

- ducotedechezswann, 3

  

- laprisonniere, 4
- lecotedeguermantes, 4
- letempretrouve, 4

  

- pr\_detect\_days, 8
- pr\_detect\_months, 9
- pr\_detect\_pro, 9
- pr\_normalize\_punc, 10
- pr\_stem\_sentences, 11
- pr\_stem\_words, 11
- proust\_books, 5
- proust\_char, 5
- proust\_characters, 6
- proust\_random, 6
- proust\_sentiments, 7
- proust\_stopwords, 7

  

- sentiments\_polarity, 12
- sentiments\_score, 12
- sodomeetgomorrhe, 13