

# Package ‘JATSdecoder’

September 27, 2022

**Title** A Metadata and Text Extraction and Manipulation Tool Set

**Date** 2022-09-26

**Version** 1.1

**Maintainer** Ingmar Böschén <ingmar.boeschén@uni-hamburg.de>

**Description** Provides a function collection to extract metadata, sectioned text and study characteristics from scientific articles in 'NISO-JATS' format. Articles in PDF format can be converted to 'NISO-JATS' with the 'Content ExtRactor and MINer' ('CER-MINE', <<https://github.com/CeON/CERMINE>>). For convenience, two functions bundle the extraction heuristics: JATSdecoder() converts 'NISO-JATS'-tagged XML files to a structured list with elements title, author, journal, history, 'DOI', abstract, sectioned text and reference list. study.character() extracts multiple study characteristics like number of included studies, statistical methods used, alpha error, power, statistical results, correction method for multiple testing, software used. An estimation of the involved sample size is performed based on reports within the abstract and the reported degrees of freedom within statistical results. In addition, the package contains some useful functions to process text (text2sentences(), text2num(), ngram(), strsplit2(), grep2()).

**Depends** R (>= 3.1.1)

**Imports** utils, stats, NLP, openNLP

**License** GPL-3

**URL** <https://github.com/ingmarboeschén/JATSdecoder>

**BugReports** <https://github.com/ingmarboeschén/JATSdecoder/issues>

**Language** en-US

**Encoding** UTF-8

**RoxygenNote** 7.2.1

**References** Böschén, I. (2021) Software review: The JATSdecoder package—extract metadata, abstract and sectioned text from NISO-JATS coded XML documents; Insights to PubMed central's open access database. Scientometrics (2021). <<https://doi.org/10.1007/s11192-021-04162-z>>

**NeedsCompilation** no

**Author** Ingmar Böschén [aut, cre] (<<https://orcid.org/0000-0003-1159-3991>>)

**R topics documented:**

allStats . . . . .	3
est.ss . . . . .	4
get.abstract . . . . .	5
get.aff . . . . .	6
get.alpha.error . . . . .	7
get.assumptions . . . . .	8
get.author . . . . .	8
get.category . . . . .	9
get.country . . . . .	10
get.doi . . . . .	10
get.editor . . . . .	11
get.history . . . . .	12
get.journal . . . . .	12
get.keywords . . . . .	13
get.method . . . . .	14
get.multi.comparison . . . . .	14
get.n.studies . . . . .	15
get.outlier.def . . . . .	16
get.power . . . . .	16
get.R.package . . . . .	17
get.references . . . . .	18
get.sig.adjectives . . . . .	18
get.software . . . . .	19
get.stats . . . . .	20
get.subject . . . . .	21
get.tables . . . . .	22
get.test.direction . . . . .	23
get.text . . . . .	23
get.title . . . . .	24
get.type . . . . .	25
get.vol . . . . .	26
grep2 . . . . .	26
has.interaction . . . . .	27
JATSdecoder . . . . .	28
letter.convert . . . . .	30
ngram . . . . .	31
standardStats . . . . .	32
strsplit2 . . . . .	33
study.character . . . . .	34
text2num . . . . .	36
text2sentences . . . . .	37
vectorize.text . . . . .	38
which.term . . . . .	39

---

`allStats`*allStats*

---

## Description

Extracts statistical results within a text string and outputs a vector of sticked results, e.g.: `c("t(12)=1.2, p>.05", "r's(33)>.7, ps<.05")`, that can be further processed with `standardStats`. This function is implemented in `get.stats` which returns the results of `allStats` and `standardStats`. Besides only plain textual input, `get.stats` enables direct processing of different file formats (NISO-JATS coded XML, DOCX, HTML) without text preprocessing.

## Usage

```
allStats(x)
```

## Arguments

`x` A character string that may contain statistical results.

## Value

Vector with sticked results. Empty, if no result is detected.

## Source

A minimal web application that extracts statistical results from single documents with `get.stats` is hosted at: <https://www.get-stats.app/>

## References

Böschchen (2021). "Evaluation of JATSdecoder as an automated text extraction tool for statistical results in scientific reports." *Scientific Reports*. doi: [10.1038/s41598-021-98782-3](https://doi.org/10.1038/s41598-021-98782-3).

## See Also

`study.character` for extracting multiple study characteristics at once.

`get.stats` for extracting statistical results from textual input and different file formats.

## Examples

```
x<-c("The mean difference of scale A was significant (beta=12.9, t(18)=2.5, p<.05)",
"The ANOVA yielded significant results on factor A (F(2,18)=6, p<.05, eta(g)2<-.22).",
"The correlation of x and y was r=.37.")
allStats(x)
```

---

*est.ss**est.ss*

---

## Description

Function to estimate studies sample size by maximizing different conservative estimates. Performs four different extraction heuristics for sample sizes mentioned in abstract, text and statistical results.

## Usage

```
est.ss(  
  abstract = NULL,  
  text = NULL,  
  stats = NULL,  
  standardStats = NULL,  
  quantileDF = 0.9,  
  max.only = FALSE,  
  max.parts = TRUE  
)
```

## Arguments

<code>abstract</code>	an abstract text string.
<code>text</code>	the main text string to process (usually method and result sections). If text has content, arguments "stats" and "standardStats" are deactivated and filled with results by <code>get.stats(text)</code> .
<code>stats</code>	statistics extracted with <code>get.stats(x)\$stats</code> (only active if no text is submitted).
<code>standardStats</code>	standard statistics extracted with <code>get.stats(x)\$standardStats</code> (only active if no text is submitted).
<code>quantileDF</code>	quantile of $(df1-1)+(df2+2)$ to extract.
<code>max.only</code>	Logical. If TRUE only the final estimate will be returned, if FALSE all sub estimates are returned as well.
<code>max.parts</code>	Logical. If FALSE outputs all captured sample sizes in sub inputs.

## Details

Sample size extraction from abstract:

- Extracts N= from abstract text and performs position-of-speech search with list of synonyms of sample units

Sample size extraction from text:

- Unifies and extracts textlines with age descriptions, than computes sum of hits as nage
- Unifies and extracts all "numeric male-female" patterns than computes sum of first male/female hit
- Unifies and extracts textlines with participant description than computes sum of first three hits as ntext

Sample size extraction from statistical results:

- Extracts "N=" in statistical results extracted with `allStats()` that contain p-value: e.g.: `chi(2, N=12)=15.2, p<.05`

Sample size extraction by degrees of freedom with result of standardStats(allStats()):  
 - Extracts df1 and df2 if possible and neither containing a ".", then calculates specified quantile of (df1+1)+(df2+2) (at least 2 group comparison assumed)

### Value

Numeric vector with extracted sample sizes by input and estimated sample size.

### See Also

[study.character](#) for extracting multiple study characteristics at once.

### Examples

```
## Not run:
a<-"One hundred twelve students participated in our study."
est.ss(abstract=a)
x<-"Our sample consists of three hundred twenty five undergraduate students.
  The F-test indicates significant differences in means F(2,102)=3.21, p<.05."
est.ss(text=x)

## End(Not run)
```

---

get.abstract

get.abstract

---

### Description

Extracts abstract tag from NISO-JATS coded XML file or text as vector of abstracts.

### Usage

```
get.abstract(
  x,
  sentences = FALSE,
  remove.title = TRUE,
  letter.convert = TRUE,
  cermine = FALSE
)
```

### Arguments

x	a NISO-JATS coded XML file or text.
sentences	Logical. If TRUE abstract is returned as vector of sentences.
remove.title	Logical. If TRUE removes section titles in abstract.
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Uni-code.
cermine	Logical. If TRUE and if 'letter.convert=TRUE' CERMINE specific letter correction is carried out (e.g. inserting of missing operators to statistical results).

**Value**

Character. The abstract/s text as floating text or vector of sentences.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <abstract>Some abstract</abstract> some text"
get.abstract(x)
x<-"Some text <abstract>Some abstract</abstract> TEXT <abstract with subsettings>
Some other abstract</abstract> Some text "
get.abstract(x)
```

---

get.aff

get.aff

---

**Description**

Extracts the affiliation tag information from NISO-JATS coded XML file or text as a vector of affiliations.

**Usage**

```
get.aff(x, remove.html = FALSE, letter.convert = TRUE)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
remove.html	Logical. If TRUE removes all html tags.
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.

**Value**

Character vector with the extracted affiliation name/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <aff>Some affiliation</aff> some text"
get.aff(x)
x<-"TEXT <aff>Some affiliation</aff> TEXT <aff>Some other affiliation</aff> TEXT"
get.aff(x)
```

---

get.alpha.error	<i>get.alpha.error</i>
-----------------	------------------------

---

## Description

Extracts reported and corrected alpha error from text and 1-alpha confidence intervals.

## Usage

```
get.alpha.error(x, p2alpha = TRUE, output = "list")
```

## Arguments

x	text string to process.
p2alpha	Logical. If TRUE detects and extracts alpha errors denoted with a critical p-value (may lead to some false positive detections).
output	One of c("list", "vector"). If output="list" returns a list containing: alpha_error, corrected_alpha, alpha_from_CI, alpha_max, alpha_min. If output="vector" returns unique alpha errors but no distinction of types.

## Value

Numeric. Vector with identified alpha-error/s.

## See Also

[study.character](#) for extracting multiple study characteristics at once.

## Examples

```
x<-c("The threshold for significance was adjusted to .05/2",  
      "Type 1 error rate was alpha=.05.")  
get.alpha.error(x)  
x<-c("We used p<.05 as level of significance.",  
      "We display .95 CIs and use an adjusted alpha of .10/3.",  
      "The effect was significant with p<.025.")  
get.alpha.error(x)
```

---

get.assumptions	<i>get.assumptions</i>
-----------------	------------------------

---

### Description

Extracts the mentioned statistical assumptions from a text string by a dictionary search of 22 common statistical assumptions.

### Usage

```
get.assumptions(x, hits_only = TRUE)
```

### Arguments

x	text string to process.
hits_only	Logical. If TRUE returns the detected assumptions only, else a hit matrix with all potential assumptions is returned.

### Value

Character. Vector with identified statistical assumption/s.

### See Also

[study.character](#) for extracting multiple study characteristics at once.

### Examples

```
x<-"Sphericity assumption and gauss-marcov was violated."
get.assumptions(x)
```

---

get.author	<i>get.author</i>
------------	-------------------

---

### Description

Extracts author tag information from NISO-JATS coded XML file or text.

### Usage

```
get.author(x, paste = "", short.names = FALSE, letter.convert = FALSE)
```



**Arguments**

<code>x</code>	a NISO-JATS coded XML file or text.
<code>paste</code>	if <code>paste!=""</code> author list is collapsed to one cell with separator specified (e.g. <code>paste=";"</code> ).
<code>short.names</code>	Logical. If TRUE fully available first names will be reduced to single letter abbreviation.
<code>letter.convert</code>	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.

**Value**

Character vector with the extracted author name/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

<code>get.category</code>	<i>get.category</i>
---------------------------	---------------------

---

**Description**

Extracts category tag/s from NISO-JATS coded XML file or text as vector of categories.

**Usage**

```
get.category(x)
```

**Arguments**

<code>x</code>	a NISO-JATS coded XML file or text.
----------------	-------------------------------------

**Value**

Character vector with the extracted category name/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <article-categories>Some category</article-categories> some text"
get.category(x)
```

---

get.country	<i>get.country</i>
-------------	--------------------

---

### Description

Extracts country tag from NISO-JATS coded XML file or text as vector of unique countries.

### Usage

```
get.country(x, unifyCountry = TRUE)
```

### Arguments

**x** a NISO-JATS coded XML file or text.  
**unifyCountry** Logical. If TRUE replaces country name with standardised country name.

### Value

Character vector with the extracted country name/s.

### See Also

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

### Examples

```
x<-"Some text <country>UK</country> some text <country>England</country>
  Text<country>Berlin, Germany</country>"
get.country(x)
```

---

get.doi	<i>get.doi</i>
---------	----------------

---

### Description

Extracts articles doi from NISO-JATS coded XML file or text.

### Usage

```
get.doi(x)
```

### Arguments

**x** a NISO-JATS coded XML file or text.

**Value**

Character string with the extracted doi.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

get.editor

*get.editor*

---

**Description**

Extracts editor tag from NISO-JATS coded XML file or text as vector of editors.

**Usage**

```
get.editor(x, role = FALSE, short.names = FALSE, letter.convert = FALSE)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
role	Logical. If TRUE adds role to editor name, if available.
short.names	Logical. If TRUE reduces fully available first names to one letter abbreviation.
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.

**Value**

Character string with the extracted editor name/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

`get.history`*get.history*

---

**Description**

Extracts available publishing history tags from NISO-JATS coded XML file or text and compute pubDate and pubyear.

**Usage**

```
get.history(x, remove.na = FALSE)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
remove.na	Logical. If TRUE hides non available date stamps.

**Value**

Character vector with the extracted dates of publishing history.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

`get.journal`*get.journal*

---

**Description**

Extracts journal tag from NISO-JATS coded XML file or text.

**Usage**

```
get.journal(x)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
---	-------------------------------------

**Value**

Character string with the extracted journal name.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <journal-title>PLoS One</journal-title> some text"
get.journal(x)
```

---

get.keywords

get.keywords

---

**Description**

Extracts keyword tag/s from NISO-JATS coded XML file or text as vector of keywords.

**Usage**

```
get.keywords(
  x,
  paste = "",
  letter.convert = TRUE,
  include.max = length(keyword)
)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
paste	if paste!="" keyword list is collapsed to one cell with separator specified (e.g. paste=";").
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.
include.max	a maximum number of keywords to extract.

**Value**

Character vector with extracted keyword/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <kwd>Keyword 1</kwd>, <kwd>Keyword 2</kwd> some text"
get.keywords(x)
get.keywords(x,paste(", "))
```

---

get.method

get.method

---

### Description

Extracts statistical methods mentioned in text.

### Usage

```
get.method(x, add = NULL, cermine = FALSE)
```

### Arguments

x	text to extract statistical methods from.
add	possible new end words of method as vector.
cermine	Logical. If TRUE CERMINE specific letter conversion will be performed.

### Value

Character. Vector with identified statistical method/s

### See Also

[study.character](#) for extracting multiple study characteristics at once.

### Examples

```
x<-"We used multiple regression analysis and  
two sample t tests to evaluate our results."  
get.method(x)
```

---

get.multi.comparison

get.multi.comparison

---

### Description

Extracts alpha-/p-value correction method for multiple comparisons from list with 15 correction methods.

### Usage

```
get.multi.comparison(x)
```

### Arguments

x	text string to process.
---	-------------------------

**Value**

Character. Identified author/method of multiple comparison correction procedure.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

**Examples**

```
x<-"We used Bonferroni corrected p-values."  
get.multi.comparison(x)
```

---

get.n.studies	<i>get.n.studies</i>
---------------	----------------------

---

**Description**

Extracts number of studies/experiments from text.

**Usage**

```
get.n.studies(x, tolower = TRUE)
```

**Arguments**

x	text string to process.
tolower	Logical. If TRUE lowerises text and search patterns for processing.

**Value**

Numeric number of identified number of studies. Returns '1' as standard output.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

---

get.outlier.def	<i>get.outlier.def</i>
-----------------	------------------------

---

### Description

Extracts outlier/extreme value definition/removal in standard deviations, if present in text.

### Usage

```
get.outlier.def(x, range = c(1, 10))
```

### Arguments

x	Character. A text string to process.
range	Numeric vector with length=2. Possible result space of extracted value/s in standard deviations. Use 'c(0,Inf)' for no restriction.

### Value

Numeric. Vector with identified outlier definition in standard deviations.

### See Also

[study.character](#) for extracting multiple study characteristics at once.

### Examples

```
x<-"We removed 4 extreme values that were 3 SD above mean."
get.outlier.def(x)
```

---

get.power	<i>get.power</i>
-----------	------------------

---

### Description

Extracts a priori power and empirical power values from text.

### Usage

```
get.power(x)
```

### Arguments

x	text string to process.
---	-------------------------



**Value**

Numeric. Identified power values.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

**Examples**

```
x<-"We used G*Power 3 to calculate the needed sample with  
beta error rate set to 12% and alpha error to .05."  
get.power(x)
```

---

get.R.package

*get.R.package*

---

**Description**

Extracts mentioned R packages from text.

**Usage**

```
get.R.package(x, update.package.list = FALSE)
```

**Arguments**

x                    text string to process.

update.package.list

Logical. If TRUE update of list with available packages is downloaded from CRAN with `utils::available.packages()`.

**Value**

Character. Vector with identified R package/s.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

**Examples**

```
get.R.package("We used the R Software packages lme4 (and psych).")
```

---

get.references	<i>get.references</i>
----------------	-----------------------

---

**Description**

Extracts reference list from NISO-JATS coded XML file or text as vector of references.

**Usage**

```
get.references(
  x,
  letter.convert = FALSE,
  remove.html = FALSE,
  extract = "full"
)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.
remove.html	Logical. If TRUE removes all HTML tags.
extract	part of references to extract (one of "full" or "title").

**Value**

Character vector with extracted references from reference list.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

get.sig.adjectives	<i>get.sig.adjectives</i>
--------------------	---------------------------

---

**Description**

Extracts adjectives used for in/significance out of list with 37 potential adjectives.

**Usage**

```
get.sig.adjectives(x, unique_only = FALSE)
```

**Arguments**

`x` text string to process.  
`unique_only` Logical. If TRUE returns unique hits only.

**Value**

Character. Vector with identified adjectives.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

**Examples**

```
get.sig.adjectives(
  x<-"We found very highly significance for type 1 effect"
)
```

---

get.software	<i>get.software</i>
--------------	---------------------

---

**Description**

Extracts mentioned software from text by dictionary search for 63 software names (object: `.software_names`).

**Usage**

```
get.software(x, add.software = NULL)
```

**Arguments**

`x` text string to process.  
`add.software` a text vector with additional software name patterns to search for.

**Value**

Character. Vector with identified statistical software/s.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

**Examples**

```
get.software("We used the R Software and Excel 4.0 to analyse our data.")
```

get.stats

*get.stats***Description**

Extracts statistical results from text string, XML, CERMXML, HTML or DOCX files. The result is a list with a vector containing all identified sticked results and a matrix containing the reported standard statistics and recalculated p-values if computation is possible.

**Usage**

```
get.stats(
  x,
  output = "both",
  stats.mode = "all",
  recalculate.p = TRUE,
  alternative = "undirected",
  estimateZ = FALSE,
  T2t = FALSE,
  R2r = FALSE,
  select = NULL,
  rm.na.col = TRUE,
  cermine = FALSE
)
```

**Arguments**

x	NISO-JATS coded XML or DOCX file path or plain textual content.
output	Select the desired output. One of c("both", "allStats", "standardStats").
stats.mode	Select a subset of test results by p-value checkability for output. One of: c("all", "checkable", "computable", "uncomputable").
recalculate.p	Logical. If TRUE recalculates p-values of test results if possible.
alternative	Character. Select sidedness of recomputed p-values from t-, r- and beta-values. One of c("undirected", "directed", "both").
estimateZ	Logical. If TRUE detected beta-/d-value is divided by reported standard error "SE" to estimate Z-value ("Zest") for observed beta/d and recompute p-value. Note: This is only valid, if Gauss-Marcov assumptions are met and a sufficiently large sample size is used. If a Z- or t-value is detected in a report of a beta-/d-coefficient with SE, no estimation will be performed, although set to TRUE.
T2t	Logical. If TRUE capital letter T is treated as t-statistic.
R2r	Logical. If TRUE capital letter R is treated as correlation.
select	Select specific standard statistics only (e.g.: c("t", "F", "Chi2")).
rm.na.col	Logical. If TRUE removes all columns with only NA from standardStats.
cermine	Logical. If TRUE CERMINE specific letter conversion will be performed on allStats results.

**Value**

If output="all": list with two elements. E1: vector of extracted results by [allStats](#) and E2: matrix of standard results by [standardStats](#).

If output="allStats": vector of extracted results by [allStats](#).

If output="standardStats": matrix of standard results by [standardStats](#).

**Source**

A minimal web application that extracts statistical results from single documents with [get.stats](#) is hosted at: <https://www.get-stats.app/>

Statistical results extracted with [get.stats](#) can be analyzed and used to identify articles stored in the PubMed Central library at: <https://www.scianalyzer.com/>.

**References**

Böschen (2021). "Evaluation of JATSdecoder as an automated text extraction tool for statistical results in scientific reports." *Scientific Reports*. doi: [10.1038/s41598-021-98782-3](https://doi.org/10.1038/s41598-021-98782-3).

**See Also**

[study.character](#) for extracting different study characteristics at once.

**Examples**

```
x<-c("The mean difference of scale A was significant (beta=12.9, t(18)=2.5, p<.05).",
"The ANOVA yielded significant results on
faktor A (F(2,18)=6, p<.05, eta(g)2<-.22)",
"the correlation of x and y was r=.37.")
get.stats(x)
```

---

get.subject

get.subject

---

**Description**

Extracts subject tag/s from NISO-JATS coded XML file or text as vector of subjects.

**Usage**

```
get.subject(x, letter.convert = TRUE, paste = "")
```

**Arguments**

x	a NISO-JATS coded XML file or text.
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Uni-code.
paste	if paste!="" subject list is collapsed to one cell with separator specified (e.g. paste=";").

**Value**

Character vector with extracted subject/s.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

**Examples**

```
x<-"Some text <subject>Some subject</subject> some text"
get.subject(x)
x<-"Some text <subject>Some subject</subject> TEXT ...
<subject>Some other subject</subject> Some text "
get.subject(x)
get.subject(x,paste=", ")
```

---

get.tables

*get.tables*

---

**Description**

Extracts HTML tables as vector of tables.

**Usage**

```
get.tables(x)
```

**Arguments**

x                      HTML file or html text.

**Value**

Character vector with extracted table in html coding.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

get.test.direction	<i>get.test.direction</i>
--------------------	---------------------------

---

**Description**

Extracts mentioned test direction/s (one sided, two sided, one and two sided) from text.

**Usage**

```
get.test.direction(x)
```

**Arguments**

x                      text string to process.

**Value**

Character.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

---

get.text	<i>get.text</i>
----------	-----------------

---

**Description**

Extracts main textual content from NISO-JATS coded XML file or text as sectioned text.

**Usage**

```
get.text(
  x,
  sectionsplit = "",
  grepsection = "",
  letter.convert = TRUE,
  greek2text = FALSE,
  sentences = FALSE,
  paragraph = FALSE,
  cermine = "auto",
  rm.table = TRUE,
  rm.formula = TRUE,
  rm.xref = TRUE,
  rm.media = TRUE,
  rm.graphic = TRUE,
  rm.ext_link = TRUE
)
```

**Arguments**

<code>x</code>	a NISO-JATS coded XML file or text.
<code>sectionsplrit</code>	search patterns for section split (forced to lower case), e.g. <code>c("intro", "method", "result", "discus")</code> .
<code>grepsection</code>	search pattern to reduce text to specific section namings only.
<code>letter.convert</code>	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.
<code>greek2text</code>	Logical. If TRUE some greek letters and special characters will be unified to textual representation (important to extract stats).
<code>sentences</code>	Logical. IF TRUE text is returned as sectioned list with sentences.
<code>paragraph</code>	Logical. IF TRUE "<New paragraph>" is added at the end of each paragraph to enable manual splitting at paragraphs.
<code>cermine</code>	Logical. If TRUE CERMINES specific error handling and letter conversion will be applied. If set to "auto" file name ending with 'cermxml\$' will set cermine=TRUE.
<code>rm.table</code>	Logical. If TRUE removes <table> tag from text.
<code>rm.formula</code>	Logical. If TRUE removes <formula> tags.
<code>rm.xref</code>	Logical. If TRUE removes <xref> tag (citing) from text.
<code>rm.media</code>	Logical. If TRUE removes <media> tag from text.
<code>rm.graphic</code>	Logical. If TRUE removes <graphic> and <fig> tag from text.
<code>rm.ext_link</code>	Logical. If TRUE removes <ext link> tag from text.

**Value**

List with two elements. 1: Character vector with section title/s, 2: Character vector with floating text of sections or list with vector of sentences per section/s if sentences=TRUE.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

<code>get.title</code>	<i>get.title</i>
------------------------	------------------

---

**Description**

Extracts article title from NISO-JATS coded XML file or text.

**Usage**

`get.title(x)`



**Arguments**

*x* a NISO-JATS coded XML file or text.

**Value**

Character string with extracted article title.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

*get.type**get.type*

---

**Description**

Extracts article type from NISO-JATS coded XML file or text.

**Usage**

```
get.type(x)
```

**Arguments**

*x* a NISO-JATS coded XML file or text.

**Value**

Character string with extracted article type.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

get.vol	<i>get.vol</i>
---------	----------------

---

**Description**

Extracts volume, first and last page from NISO-JATS coded XML file or text.

**Usage**

```
get.vol(x)
```

**Arguments**

x	a NISO-JATS XML coded file or text.
---	-------------------------------------

**Value**

Character string with extracted journal volume.

**See Also**

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

---

grep2	<i>grep2</i>
-------	--------------

---

**Description**

Extension of `grep()`. Allows to identify and extract cells with/without multiple search patterns that are connected with AND.

**Usage**

```
grep2(pattern, x, value = TRUE, invert = FALSE, perl = FALSE)
```

**Arguments**

pattern	Character vector containing regular expression as cells to be matched in the given character vector.
x	A character vector where matches are sought, or an object which can be coerced by <code>as.character</code> to a character vector. Long vectors are supported.
value	Logical. If FALSE, a vector containing the (integer) indices of the matches determined by <code>grep2</code> is returned, and if TRUE, a vector containing the matching elements themselves is returned.
invert	Logical. If TRUE return indices or values for elements that do not match.
perl	Logical. Should Perl-compatible regexps be used?

**Value**

grep2(value = FALSE) returns a vector of the indices of the elements of x that yielded a match (or not, for invert = TRUE). This will be an integer vector unless the input is a long vector, when it will be a double vector.

grep2(value = TRUE) returns a character vector containing the selected elements of x (after coercion, preserving names but no other attributes).

**See Also**

[grep](#)

**Examples**

```
x<-c("ab", "ac", "ad", "bc", "bad")
grep2(c("a", "b"), x)
grep2(c("a", "b"), x, invert=TRUE)
grep2(c("a", "b"), x, value=FALSE)
```

---

has.interaction	<i>has.interaction</i>
-----------------	------------------------

---

**Description**

Identifies mentions of interaction/moderator/mediator effect in text.

**Usage**

```
has.interaction(x)
```

**Arguments**

x                      text string to process.

**Value**

Character vector with type/s of identified interaction/moderator/mediator effect.

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

JATSdecoder

*JATSdecoder***Description**

Function to extract and restructure NISO-JATS coded XML file or text into a list with metadata and text as selectable elements. Use **CERMINER** to convert PDF to CERMLXML files.

**Usage**

```
JATSdecoder(
  x,
  sectionsplit = c("intro", "method", "result", "study", "experiment", "conclu",
    "implica", "discussion"),
  grepsection = "",
  sentences = FALSE,
  paragraph = FALSE,
  abstract2sentences = TRUE,
  output = "all",
  letter.convert = TRUE,
  unify.country.name = TRUE,
  greek2text = FALSE,
  warning = TRUE,
  countryconnection = FALSE,
  authorconnection = FALSE
)
```

**Arguments**

x	a NISO-JATS coded XML file or text.
sectionsplit	search patterns for section split of text parts (forced to lower case), e.g. c("intro", "method", "result", "discus").
grepsection	search pattern in regex to reduce text to specific section only.
sentences	Logical. IF TRUE text is returned as sectioned list with sentences.
paragraph	Logical. IF TRUE "<New paragraph>" is added at the end of each paragraph to enable manual splitting at paragraphs.
abstract2sentences	Logical. IF TRUE abstract is returned as vector with sentences.
output	selection of specific results to output c("all", "title", "author", "affiliation", "journal", "volume", "editor", "doi", "type", "history", "country", "subject", "keywords", "abstract", "sections", "text", "tables", "captions", "references").
letter.convert	Logical. If TRUE converts hexadecimal and HTML coded characters to Unicode.
unify.country.name	Logical. If TRUE tries to unify country name/s with list of country names from worldmap().

<code>greek2text</code>	Logical. If TRUE converts and unifies several greek letters to textual representation, e.g.: "alpha".
<code>warning</code>	Logical. If TRUE outputs a warning if processing CERMINE converted PDF files.
<code>countryconnection</code>	Logical. If TRUE outputs country connections as vector <code>c("A - B", "A - C", ...)</code> .
<code>authorconnection</code>	Logical. If TRUE outputs connections of a maximum of 50 involved authors as vector <code>c("A - B", "A - C", ...)</code> .

**Value**

List with extracted meta data, sectioned text and references.

**Note**

A short tutorial on how to work with JATSdecoder and the generated outputs can be found at: <https://github.com/ingmarboesch/JATSdecoder>

**Source**

An interactive web application for selecting and analyzing extracted article metadata and study characteristics for articles linked to PubMed Central is hosted at: <https://www.scianalyzer.com/>

The XML version of PubMed Central database articles can be downloaded in bulk from: [https://ftp.ncbi.nlm.nih.gov/pub/pmc/oa\\_bulk/](https://ftp.ncbi.nlm.nih.gov/pub/pmc/oa_bulk/)

**References**

Böschchen (2021). "Software review: The JATSdecoder package - extract metadata, abstract and sectioned text from NISO-JATS coded XML documents; Insights to PubMed Central's open access database." *Scientometrics*. doi: [10.1007/s1119202104162z](https://doi.org/10.1007/s1119202104162z).

**See Also**

[study.character](#) for extracting different study characteristics at once.

[get.stats](#) for extracting statistical results from textual input and different file formats.

**Examples**

```
## Not run:
# download example XML file via URL
x<-"https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0114876&type=manuscript"
download.file(x,"file.xml")
# convert full article to list with metadata, sectioned text and reference list
JATSdecoder("file.xml")
# extract specific content (here: abstract and text)
JATSdecoder("file.xml",output=c("abstract","text"))
# or use specific functions, e.g.:
get.abstract("file.xml")
```

```
get.text("file.xml")  
  
## End(Not run)
```

---

letter.convert	<i>letter.convert</i>
----------------	-----------------------

---

**Description**

Converts and unifies most hexadecimal and some HTML coded letters to Unicode characters. Performs CERMINE specific error correction (inserting operators, where these got lost while conversion).

**Usage**

```
letter.convert(x, cermine = FALSE, greek2text = FALSE, warning = TRUE)
```

**Arguments**

x	text string to process.
cermine	Logical. If TRUE CERMINE specific error handling and letter conversion will be applied.
greek2text	Logical. If TRUE some greek letters and special characters will be unified to textual representation (important to extract stats).
warning	Logical. If TRUE prints warning message if CERMINE specific letter conversion was performed.

**Value**

Character. Text with unified and corrected letter representation.

**Examples**

```
x<-c("five &#x0003c; ten","five &lt; ten")  
letter.convert(x)
```

ngram

*ngram***Description**

Extracts ngram bag of words around words that match a search pattern. Note: If an input contains the search pattern twice, only the ngram bag of words of the last hit is detected. Consider individual text splitting with `text2sentences()` or `strsplit2()` before applying `ngram()`.

**Usage**

```
ngram(
  x,
  pattern,
  ngram = c(-3, 3),
  tolower = FALSE,
  split = FALSE,
  exact = FALSE
)
```

**Arguments**

<code>x</code>	vector of text strings to process.
<code>pattern</code>	a search term pattern to extract the ngram bag of words.
<code>ngram</code>	a vector of length=2 that defines the number of words to extract from left and right side of pattern match.
<code>tolower</code>	Logical. If TRUE converts text and pattern to lower case.
<code>split</code>	Logical. If TRUE splits text input at "[.,;:] " before processing. Note: You may consider other text splits before.
<code>exact</code>	Logical. If TRUE only exact word matches will be proceses

**Value**

Character. Vector with +-n words of search pattern.

**Examples**

```
text<-"One hundred twenty-eight students participated in our Study,
that was administred in thirteen clinics."
ngram(text,pattern="study",ngram=c(-1,2))
```

---

standardStats

*standardStats*


---

## Description

Extracts and restructures statistical standard results like Z, t, Cohen's d, F, eta<sup>2</sup>, r, R<sup>2</sup>, chi<sup>2</sup>, BF<sub>10</sub>, Q, U, H, OR, RR, beta values into a matrix. Performs a recomputation of two- and one-sided p-values if possible. This function is implemented in [get.stats](#) which returns the results of [allStats](#) and [standardStats](#). Besides only plain textual input, [get.stats](#) enables direct processing of different file formats (NISO-JATS coded XML, DOCX, HTML) without text preprocessing.

## Usage

```
standardStats(
  x,
  stats.mode = "all",
  recalculate.p = TRUE,
  alternative = "undirected",
  estimateZ = FALSE,
  T2t = FALSE,
  R2r = FALSE,
  select = NULL,
  rm.na.col = TRUE
)
```

## Arguments

x	result vector by <a href="#">allStats</a> or chracter vector with a statistical test result per cell, e.g. c("t(12)=1.2, p>.05", "chi2(2)=12.7, p<.05")
stats.mode	Select subset of standard stats. One of: c("all", "checkable", "computable", "uncomputable").
recalculate.p	Logical. If TRUE recalculates p values (for 2 sided test) if possible.
alternative	Character. Select sidedness of recomputed p-values from t-, r- and beta-values. One of c("undirected", "directed", "both").
estimateZ	Logical. If TRUE detected beta-/d-value is divided by reported standard error "SE" to estimate Z-value ("Zest") for observed beta/d and recompute p-value. Note: This is only valid, if Gauss-Marcov assumptions are met and a sufficiently large sample size is used. If a Z- or t-value is detected in a report of a beta-/d-coefficient with SE, no estimation will be performed, although set to TRUE.
T2t	Logical. If TRUE capital letter T is treated as t-statistic.
R2r	Logical. If TRUE capital letter R is treated as correlation.
select	Select specific standard statistics only (e.g.: c("t", "F", "Chi2")).
rm.na.col	Logical. If TRUE removes all columns with only NA.



**Value**

Matrix with recognized statistical standard results and recalculated p-values. Empty, if no result is detected.

**Source**

A minimal web application that extracts statistical results from single documents with [get.stats](#) is hosted at: <https://www.get-stats.app/>

Statistical results extracted with [get.stats](#) can be analyzed and used to identify articles stored in the PubMed Central library at: <https://www.scianalyzer.com/>.

**References**

Böschchen (2021). "Evaluation of JATSdecoder as an automated text extraction tool for statistical results in scientific reports." *Scientific Reports*. doi: [10.1038/s41598-021-98782-3](https://doi.org/10.1038/s41598-021-98782-3).

**See Also**

[study.character](#) for extracting multiple study characteristics at once.

[get.stats](#) for extracting statistical results from textual input and different file formats.

**Examples**

```
x<-c("t(38.8)<=>1.96, p<=>.002", "F(2,39)<=>4, p<=>.05",  
      "U(2)=200, p>.25", "Z=2.1, F(20.8,22.6)=200, p<.005,  
      BF(01)>4", "chi=3.2, r(34)=-.7, p<.01, R2=76%.")  
standardStats(x)
```

---

strsplit2

*strsplit2*

---

**Description**

Extension of strsplit(). Makes it possible to split lines before or after a pattern match without removing the pattern.

**Usage**

```
strsplit2(x, split, type = "remove", perl = FALSE)
```

**Arguments**

x	text string to process.
split	pattern to split text at.
type	one out of c("remove", "before", "after").
perl	Logical. If TRUE uses perl expressions.

Value

A list of the same length as x, the i-th element of which contains the vector of splits of x[i].

Examples

```
x<-"This is some text, where text is the split pattern of the text."
strsplit2(x,"text","after")
```

---

study.character	<i>study.character</i>
-----------------	------------------------

---

Description

Extracts study characteristics out of a NISO-JATS coded XML file. Use **CERMINE** to convert PDF to CERMLXML files.

Usage

```
study.character(  
  x,  
  stats.mode = "all",  
  recalculate.p = TRUE,  
  alternative = "auto",  
  estimateZ = FALSE,  
  T2t = FALSE,  
  R2r = FALSE,  
  selectStandardStats = NULL,  
  p2alpha = TRUE,  
  alpha_output = "list",  
  captions = TRUE,  
  text.mode = 1,  
  update.package.list = FALSE,  
  add.software = NULL,  
  quantileDF = 0.9,  
  N.max.only = FALSE,  
  output = "all",  
  rm.na.col = TRUE  
)
```

Arguments

x	NISO-JATS coded XML file.
stats.mode	Character. Select subset of standard stats. One of: c("all", "checkable", "computable").
recalculate.p	Logical. If TRUE recalculates p values (for 2 sided test) if possible.

alternative	Character. Select sidedness of recomputed p-values for t-, r- and Z-values. One of c("auto", "undirected", "directed", "both"). If set to "auto" 'alternative' will be set to 'both' if get.test.direction() detects one-directional hypotheses/tests in text. If no directional hypotheses/tests are detected only "undirected" recomputed p-values will be returned.
estimateZ	Logical. If TRUE detected beta-/d-value is divided by reported standard error "SE" to estimate Z-value ("Zest") for observed beta/d and recompute p-value. Note: This is only valid, if Gauss-Marcov assumptions are met and a sufficiently large sample size is used. If a Z- or t-value is detected in a report of a beta-/d-coefficient with SE, no estimation will be performed, although set to TRUE.
T2t	Logical. If TRUE capital letter T is treated as t-statistic when extracting statistics with get.stats().
R2r	Logical. If TRUE capital letter R is treated as correlation when extracting statistics with get.stats().
selectStandardStats	Select specific standard statistics only (e.g.: c("t", "F", "Chi2")).
p2alpha	Logical. If TRUE detects and extracts alpha errors denoted with critical p-value (what may lead to some false positive detections).
alpha_output	One of c("list", "vector"). If alpha_output = "list" a list with elements: alpha_error, corrected_alpha, alpha_from_CI, alpha_max, alpha_min is returned. If alpha_output = "vector" unique alpha errors without a distinction of types is returned.
captions	Logical. If TRUE captions text will be scanned for statistical results.
text.mode	Numeric. Defines text parts to extract statistical results from (text.mode=1: abstract and full text, text.mode=2: method and result section, text.mode=3: result section only).
update.package.list	Logical. If TRUE updates available R packages with utils::available.packages() function.
add.software	additional software names to detect as vector.
quantileDF	quantile of (df1+1)+(df2+1) to extract for estimating sample size.
N.max.only	return only maximum of estimated sample sizes.
output	output selection of specific results c("doi", "title", "year", "Nstudies", "methods", "alpha_error", "power", "multi_comparison_correction", "assumptions", "OutlierRemovalInSD", "InteractionModeratorMediatorEffect", "test_direction", "sig_adjectives", "software", "Rpackage", "stats", "standardStats", "estimated_sample_size").
rm.na.col	Logical. If TRUE removes all columns with only NA in extracted standard statistics.

## Value

List with extracted study characteristics.

## Note

A short tutorial on how to work with JATSdecoder and the generated outputs can be found at: <https://github.com/ingmarboesch/JATSdecoder>

## Source

An interactive web application for selecting and analyzing extracted article metadata and study characteristics for articles linked to PubMed Central is hosted at: <https://www.scianalyzer.com/>

The XML version of PubMed Central database articles can be downloaded in bulk from: [https://ftp.ncbi.nlm.nih.gov/pub/pmc/oa\\_bulk/](https://ftp.ncbi.nlm.nih.gov/pub/pmc/oa_bulk/)

## References

Böschchen (2021). "Evaluation of JATSdecoder as an automated text extraction tool for statistical results in scientific reports." *Scientific Reports*. doi: [10.1038/s41598-021-98782-3](https://doi.org/10.1038/s41598-021-98782-3).

## See Also

[JATSdecoder](#) for simultaneous extraction of meta-tags, abstract, sectioned text and reference list.

[get.stats](#) for extracting statistical results from textual input and different file formats.

## Examples

```
## Not run:
# download example XML file via URL
x<-"https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0114876&type=manuscript"
download.file(x,"file.xml")
# convert full article to list with study characteristics
study.character("file.xml")

## End(Not run)
```

---

text2num

text2num

---

## Description

Converts special annotated number and written numbers in a text string to a fully digit representation. Can handle numbers with exponent, fraction, percent, e+num, products and written representation (e.g. 'fourty-one') of all absolute numbers up to 99,999 (Note: gives wrong output for higher spelled numbers). Process is performed in the same order as its arguments.

**Usage**

```
text2num(
  x,
  exponent = TRUE,
  percentage = TRUE,
  fraction = TRUE,
  e = TRUE,
  product = TRUE,
  words = TRUE
)
```

**Arguments**

x	text string to process.
exponent	Logical. If TRUE values with exponent are converted to a digit representation.
percentage	Logical. If TRUE percentages are converted to a digit representation.
fraction	Logical. If TRUE fractions are converted to a digit representation.
e	Logical. If TRUE values denoted with 'number e+number' (e.g. '2e+2') or number*10^number are converted to a digit representation.
product	Logical. If TRUE values products are converted to a digit representation.
words	Logical. If TRUE written numbers are converted to a digit representation.

**Value**

Character. Text with unified digital representation of numbers.

**Examples**

```
x<-c("numbers with exponent: 2^2, -2.5^2, (-3)^2, 6.25^.5, .2^-2 text.",
      "numbers with percentage: 2%, 15 %, 25 percent.",
      "numbers with fractions: 1/100, -2/5, -7/.1",
      "numbers with e: 10e+2, -20e3, .2E-2, 2e4",
      "numbers as products: 100*2, -20*.1, 2*10^3",
      "written numbers: twenty-two, one hundred forty five, fifteen percent",
      "mix: one hundred ten is not 1/10 is not 10^2 nor 10%/5")
text2num(x)
```

---

text2sentences

text2sentences

---

**Description**

Converts floating text to a vector with sentences via fine-tuned regular expressions.

**Usage**

```
text2sentences(x)
```

**Arguments**

x                      text string to process.

**Value**

Character vector with sentences compiled from floating text.

**Examples**

```
x<-"Some text with result (t(18)=1.2, p<.05). This shows how text2sentences works."  
text2sentences(x)
```

---

vectorize.text

*vectorize.text*

---

**Description**

Converts vector of text to a list of vectors with words within each cell. Note: punctuation will be removed.

**Usage**

```
vectorize.text(x)
```

**Arguments**

x                      text string to vectorize.

**Value**

Character vector with one word per cell.

**Examples**

```
text<-"One hundred twenty-eight students participated in our  
Study, that was administred in thirteen clinics."  
vectorize.text(text)
```

---

<code>which.term</code>	<i><code>which.term</code></i>
-------------------------	--------------------------------

---

**Description**

Returns search element/s from vector that is/are present in text or returns search term hit vector for all terms.

**Usage**

```
which.term(x, terms, tolower = TRUE, hits_only = FALSE)
```

**Arguments**

<code>x</code>	text string to process.
<code>terms</code>	search term vector.
<code>tolower</code>	Logical. If TRUE converts search terms and text to lower case.
<code>hits_only</code>	Logical. If TRUE returns search pattern/s, that were found in text and not a search term hit vector.

**Value**

Binary hit vector with search term named elements if `hits_only=FALSE`.  
Character vector with identified search term elements if `hits_only=TRUE`.

**Examples**

```
text<-c("This demo demonstrates how which.term works.",
        "The result is a simple 0, 1 coded vector for all search patterns or
        a vector including the identified patterns only.")
which.term(text,c("Demo","example","work"))
which.term(text,c("Demo","example","work"),tolower=TRUE,hits_only=TRUE)
```

# Index

allStats, [3](#), [3](#), [21](#), [32](#)

est.ss, [4](#)

get.abstract, [5](#)

get.aff, [6](#)

get.alpha.error, [7](#)

get.assumptions, [8](#)

get.author, [8](#)

get.category, [9](#)

get.country, [10](#)

get.doi, [10](#)

get.editor, [11](#)

get.history, [12](#)

get.journal, [12](#)

get.keywords, [13](#)

get.method, [14](#)

get.multi.comparison, [14](#)

get.n.studies, [15](#)

get.outlier.def, [16](#)

get.power, [16](#)

get.R.package, [17](#)

get.references, [18](#)

get.sig.adjectives, [18](#)

get.software, [19](#)

get.stats, [3](#), [20](#), [21](#), [29](#), [32](#), [33](#), [36](#)

get.subject, [21](#)

get.tables, [22](#)

get.test.direction, [23](#)

get.text, [23](#)

get.title, [24](#)

get.type, [25](#)

get.vol, [26](#)

grep, [27](#)

grep2, [26](#)

has.interaction, [27](#)

JATSdecoder, [6](#), [9–13](#), [18](#), [22](#), [24–26](#), [28](#), [36](#)

letter.convert, [30](#)

ngram, [31](#)

standardStats, [3](#), [21](#), [32](#), [32](#)

strsplit2, [33](#)

study.character, [3](#), [5](#), [7](#), [8](#), [14–17](#), [19](#), [21](#),  
[23](#), [27](#), [29](#), [33](#), [34](#)

text2num, [36](#)

text2sentences, [37](#)

vectorize.text, [38](#)

which.term, [39](#)