

# Package: kanjistat (via r-universe)

September 7, 2024

**Type** Package

**Title** A Statistical Framework for the Analysis of Japanese Kanji Characters

**Version** 0.14.1.9000

**Date** 2024-06-04

**Maintainer** Dominic Schuhmacher

`<dominic.schuhmacher@mathematik.uni-goettingen.de>`

**Description** Various tools and data sets that support the study of kanji, including their morphology, decomposition and concepts of distance and similarity between them.

**URL** <https://dschuhmacher.github.io/kanjistat/>

**BugReports** <https://github.com/dschuhmacher/kanjistat/issues>

**Depends** R ( $\geq 4.1$ )

**Imports** methods, graphics, grDevices, utils, crayon, dendextend, gsubfn, Matrix, png, purrr, RANN, rlang, ROI, sysfonts, showtext, stringi, stringr, transport ( $\geq 0.15$ ), xml2, lifecycle, Rcpp

**Suggests** dplyr, jsonlite, kanjistat.data, knitr, rmarkdown, ROI.plugin.glpk, systemfonts, testthat ( $\geq 3.0.0$ ), tibble, withr

**Additional\_repositories** <https://dschuhmacher.github.io/drat>

**License** GPL ( $\geq 3$ )

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 7.3.1

**VignetteBuilder** knitr

**Config/testthat/edition** 3

**Roxygen** list(markdown = TRUE)

**LinkingTo** Rcpp

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

**RemoteUrl** <https://github.com/dschuhmacher/kanjimat>

**RemoteRef** HEAD

**RemoteSha** cdf58de0ddf900fe239921c5f48fb33e2622d47d

## Contents

cjk_escape	2
codepoint	3
compare_neighborhoods	4
convert_kanji	5
distdata	6
fivebetas	7
fivetrees	8
get_strokes	9
get_strokes_compo	9
kanjidata	10
kanjidist	12
kanjidistmat	14
kanjimat	15
kanjivec	17
kmatdist	20
kmatdistmat	21
kreadmean	22
lookup	22
options	23
plot.kanjimat	24
plot.kanjivec	25
plotkanji	26
pooled_similarity	28
print.kanjivec	29
read_kanjidic2	29
samplekan	31
sedist	32
str.kanjivec	33
<b>Index</b>	<b>34</b>

---

cjk\_escape

*Replace CJK characters in files by escape sequences*

---

### Description

All CJK characters in the file(s) found at the specified path are substituted by their Unicode escape sequences (`\u` + 4 digit hex number or `\U` + 8 digit hex number where necessary).

**Usage**

```
ckj_escape(path, outdir = NULL, verbose = TRUE)
```

**Arguments**

<b>path</b>	the path to a directory or a single file.
<b>outdir</b>	the directory where the output files are written. Defaults to the subdirectory out of the directory in <b>path</b> . The output files have the same names as the originals.
<b>verbose</b>	whether to print a message for each output file.

**Details**

If **path** is a directory, the replacement is performed for all files at that location (subdirectories are ignored). If **outdir** is the same as **path**, the original files are overwritten without warning.

If **path** is a file, the replacement is limited to this file. If **outdir** is the same as `dirname(path)`, the files are overwritten without warning.

**Value**

No return value, called for side effects.

---

**codepoint**
*Convert between Unicode codepoint and kanji*


---

**Description**

Given codepoints **cp**, the function `codepointToKanji` transforms to UTF-8, which will typically show as the actual character the codepoints stands for. Vice versa, given (UTF-8 encoded) kanjis **kan**, the function `kanjiToCodepoint` transforms to unicode codepoints.

**Usage**

```
codepointToKanji(cp, concat = FALSE)
```

```
kanjiToCodepoint(kan, character = FALSE)
```

**Arguments**

<b>cp</b>	a vector of character strings or objects of class <code>hexmode</code> , representing hexadecimal numbers.
<b>concat</b>	logical. Shall the returned characters be concatenated?
<b>kan</b>	a vector of kanjis (strings of length 1) or a single string of length $\geq 1$ of kanjis.
<b>character</b>	logical. Shall the returned codepoints be of class "character" or <code>hexmode</code> .

**Value**

For `codepointToKanji` a character vector of kanji. For `kanjiToCodepoint` a vector of hexadecimal numbers (class `hexmode`).

**Examples**

```
codepointToKanji(c("51b7", "6696", "71b1"))
kanjiToCodepoint(" ")
```

---

`compare_neighborhoods`

*Compare distances of nearest kanji*

---

**Description**

List distances to nearest neighbors of a given kanji in terms of a reference distance (which is currently only the stroke edit distance) and compare with values in terms of another distance (currently only the component transport distance, a.k.a. kanji distance).

**Usage**

```
compare_neighborhoods(
  kan,
  refdist = "strokedist",
  refnn = 10,
  compdist = "kanjidist",
  compnn = 0,
  ...
)
```

**Arguments**

<code>kan</code>	a kanji (currently only as a single UTF-8 character).
<code>refdist</code>	the name of the reference distance (currently only "strokedist").
<code>refnn</code>	the number of nearest neighbors in terms of the reference distance.
<code>compdist</code>	a character vector. The name(s) of one or several other distances to compare with (currently only "kanjidist").
<code>compnn</code>	the number of nearest neighbors in terms of the other distance(s). If this is positive it is assumed that the suggested package <code>kanjistat.data</code> is available.
<code>...</code>	further parameters that are passed to <code>kanjidist()</code> .

**Value**

A matrix of distances with `refnn + compnn` columns named by the nearest neighbors of `kan` (first in terms of the reference distance, then the other distances) and `1 + length(compdist)` rows named by the type of distance.

**Warning****[Experimental]**

This is only a first draft of the function and its interface and details may change considerably in the future. As there is currently no precomputed kanjidist matrix, there is a huge difference in computation time between setting `compnn = 0` (only kanji distances to the `refnn` nearest neighbors in terms of `refdist` have to be computed) and setting `compnn` to any value  $> 0$  (kanji distances to all 2135 other Jouyou kanji have to be computed in order to determine the `compnn` nearest neighbors; depending on the system and parameter settings this can take (roughly) anywhere between 2 minutes and an hour).

**Examples**

```
# compare_neighborhoods(" ", refnn=5, compo_seg_depth=4, approx="pcweighted",
#                       compnn=0, minor_warnings=FALSE)
```

---

<code>convert_kanji</code>	<i>Convert between kanji formats</i>
----------------------------	--------------------------------------

---

**Description**

Accept any interpretable representation of kanji in terms of index numbers, UTF-8 character strings of length 1, UTF-8 codepoints or `kanjivec` objects and convert it to all or any of these formats.

**Usage**

```
convert_kanji(
  key,
  output = c("all", "index", "character", "hexmode", "kanjivec"),
  simplify = TRUE
)
```

**Arguments**

<code>key</code>	an atomic vector or list of kanji in any combination of formats.
<code>output</code>	a string describing the desired output.
<code>simplify</code>	logical. Whether to simplify the output to an atomic vector or keep the structure of the original vector. In either case it depends on output whether this is possible.

**Details**

Index numbers are in terms of the order in `kbase`. UTF-8 codepoints are usually of class "hexmode", but character strings starting with "0x" or "0X" are also accepted in the `key`.

For `output = "kanjivec"`, the GitHub package `kanjistat.data` has to be available or an error is returned. For `output = "all"`, component `kanjivec` is set to `NA` if `kanjistat.data` is not available.

**Value**

A vector of the same length as key. If `simplify` is `TRUE`, this is an atomic vector for output = "index", "character" or "hexmode", and a list for output = "kanjivec" or "all" a list. If `simplify` is `FALSE`, the original structure (atomic or list) kept whenever possible.

**Examples**

```
convert_kanji(as.hexmode("99ac"))
convert_kanji("0x99ac") # same
convert_kanji(500, "character") == kbase$kanji[500] # TRUE
```

---

<code>distdata</code>	<i>Precomputed kanji distances</i>
-----------------------	------------------------------------

---

**Description**

Precomputed kanji distances

**Usage**

`dstrokedist`

`dyehli`

**Format**

Symmetric sparse matrices containing distances between a key kanji, its ten nearest neighbors and possibly some other close kanji. For `dstrokedist`, these are the stroke edit distances according to Yencken and Baldwin (2008). For `dyehli`, these are the bag-of-radicals distances according to Yeh and Li (2002). Both are an instance of the S4 class `dsCMatrix` (symmetric sparse matrices in *column*-compressed format) with 2133 rows and 2133 columns.

All pre-2010 jouyou kanji that are also post-2010 jouyou kanji are included. The indices are those from `kbase`.

**Source**

Datasets from <https://lars.yencken.org/datasets>, made available under the Creative Commons Attribution 3.0 Unported licence.

Computed as part of Yencken, Lars (2010) *Orthographic support for passing the reading hurdle in Japanese*. PhD Thesis, University of Melbourne, Melbourne, Australia.

**References**

Yeh, Su-Ling and Li, Jing-Ling (2002). Role of structure and component in judgements of visual similarity of Chinese characters. *Journal of Experimental Psychology: Human Perception and Performance*, **28**(4), 933–947.

Yencken, Lars, & Baldwin, Timothy (2008). Measuring and predicting orthographic associations: Modelling the similarity of Japanese kanji. In: *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pp. 1041-1048.

## Examples

```
# Find index for kanji
bu_index <- match(" ", kbase$kanji)

# Look up available stroke edit distances for .
non_zero <- which(dstrokedist[bu_index,] != 0)
sed <- dstrokedist[non_zero, bu_index]
names(sed) <- kbase[non_zero,]$kanji
sort(sed)

# Look up available bag-of-radicals distances for .
non_zero <- which(dyehli[bu_index,] != 0)
bord <- dyehli[non_zero, bu_index]
names(bord) <- kbase[non_zero,]$kanji
sort(bord)
```

---

fivebetas

*A sample list of kanjivec objects*

---

## Description

A sample list of kanjivec objects

## Usage

```
fivebetas
```

## Format

`fivebetas` is a list of five `kanjivec` objects representing the basic kanji 一, 丨, 丿, 乙, 廴, containing "beta" components, which come in fact from two different classical radicals:

- 一 on the left: mound, small village
- 丨 on the right: large village

## Source

The list has been generated with the function `kanjivec` with parameter `flatten="intelligent"` from the corresponding files in the KanjiVG database by Ulrich Apel (<https://kanjivg.tagaini.net/>).

## Examples

```
oldpar <- par(mfrow = c(1,5), mai = rep(0,4))
invisible( lapply(fivebetas, plot, seg_depth = 2) )
par(oldpar)
```

---

**fivetrees***Sample lists of kanjimat objects*

---

## Description

Sample lists of kanjimat objects

## Usage

```
fivetrees1
```

```
fivetrees2
```

```
fivetrees3
```

## Format

`fivetrees1`, `fivetrees2` and `fivetrees3` are lists of five [kanjimat](#) objects each, representing the same five basic kanji 一, 二, 三, 四, 五, containing each a tree component. Their matrices are antialiased 64 x 64 pixel representations of the kanji. The size is chosen as a compromise between aesthetics and memory/computational cost, such as for [kmatdist](#).

All of them are in handwriting style fonts. `fivetrees1` is in a Kyoukasho font (schoolbook style), `fivetrees2` is in a Kaisho font (regular script calligraphy font), `fivetrees3` is in a Gyousho font (semi-cursive calligraphy font).

An object of class `list` of length 5.

An object of class `list` of length 5.

An object of class `list` of length 5.

## Source

The list has been generated with the function [kanjimat](#) using the Mac OS pre-installed YuKyokasho font (`fivetrees1`), as well as the freely available fonts `nagayama_kai` by Norio Nagayama and `KouzanBrushFontGyousyo` by Aoyagi Kouzan.

## Examples

```
oldpar <- par(mfrow = c(3,5))
invisible( lapply(fivetrees1, plot) )
invisible( lapply(fivetrees2, plot) )
invisible( lapply(fivetrees3, plot) )
par(oldpar)
```



---

get_strokes	<i>Get the strokes of a kanjivec object</i>
-------------	---

---

### Description

The strokes are the leaves of the kanjivec `stroketre`. They consist of a two-column matrix giving a discretized path for the stroke in the unit square  $[0, 1]^2$  with further attributes.

### Usage

```
get_strokes(kvec, which = 1:kvec$nstrokes, simplify = TRUE)
```

### Arguments

<code>kvec</code>	an object of class <code>kanjivec</code>
<code>which</code>	a numeric vector specifying the numbers of the strokes that are to be returned. Defaults to all strokes.
<code>simplify</code>	logical. Shall only the stroke be returned if <code>which</code> has length 1?

### Value

Usually a list of strokes with attributes. Regardless of whether `which` is ordered or contains duplicates, the returned list will always contain the strokes in their natural order without duplicates. If `which` has length 1 and `simplified = TRUE`, the list is avoided, and only the single stroke is returned.

### See Also

[get\\_strokes\\_compo](#)

### Examples

```
kanji <- fivebetas[[5]]
get_strokes(kanji, c(3,10)) # the two long vertical strokes in
```

---

get_strokes_compo	<i>Get the strokes of a specific component of a kanjivec object</i>
-------------------	---

---

### Description

The strokes are the leaves of the kanjivec `stroketre`. They consist of a two-column matrix giving a discretized path for the stroke in the unit square  $[0, 1]^2$  with further attributes.

### Usage

```
get_strokes_compo(kvec, which = c(1, 1))
```

**Arguments**

**kvec** an object of class `kanjivec`

**which** a vector of length 2 specifying the index of the component, i.e. the component used is `pluck(kvec$components, !!!which)`. The default `c(1,1)` refers to the root component (full kanji), so all strokes are returned.

**Value**

A list of strokes with attributes.

**See Also**

[get\\_strokes](#)

**Examples**

```
kanji <- fivebetas[[5]]
# get the three strokes of the component in
rad <- get_strokes_compo(kanji, c(2,1))
plot(0.5, 0.5, xlim=c(0,1), ylim=c(0,1), type="n", asp=1, xaxs="i", yaxs="i", xlab="", ylab="")
invisible(lapply(rad, lines, lwd=4))
```

---

kanjidata

*Data on kanji*

---

**Description**

The tibbles `kbase` and `kmorph` provide basic and morphologic information, respectively, for all kanji contained in the KANJIDIC2 file (see below)

**Usage**

`kbase`

`kmorph`

**Format**

`kbase` is a tibble with 13,108 rows and 13 variables:

**kanji** the kanji

**unicode** the Unicode codepoint

**strokes** the number of strokes

**class** one of four classes: "kyouiku", "jouyou", "jinmeiyou" or "hyougai"

**grade** a number from 1-11, basically a finer version of class, same as in KANJIDIC2, except that we assigned an 11 for all hyougaiji (rather than an NA value)

**kanken** at what level the kanji appears in the Nihon Kanji Nouryoku Kentei (Kanken)

**jlpt** at what level the kanji appears in the Japanese Language Proficiency Test (Nihongou Nouryoku Shiken)

**wanikani** at what level the kanji is learned on the kanji learning website Wanikani

**frank** the frequency rank (1 = most frequent) "based on several averages (Wikipedia, novels, newspapers, ...)"

**frank\_news** the frequency rank (1 = most frequent) based on news paper data (2501 most frequent kanji over four years in the Mainichi Shimbun)

**read\_on, read\_kun** a single ON reading in katakana

**read\_kun** a single kun reading in hiragana

**mean** a single English meaning of the kanji

kmorph is a tibble with 13,108 rows and 15 variables:

**kanji** the kanji

**strokes** the number of strokes

**radical** the traditional (Kangxi) radical used for indexing kanji (one of 214)

**radvar** the variant of the radical if it is different, otherwise NA

**nelson\_c** the Nelson radical if it differs from the traditional one, otherwise NA

**idc** ideographic description character (plus sometimes a number or a letter) describing the shape of the kanji

**components** visible components of the kanji; originally from KRADFILE

**skip** the kanji's SKIP code

**mean** a single English meaning of the kanji (same as in kbase)

## Details

The single ON and kun readings and the single meaning are for easy identification of the more difficult kanji. They are the first entry in the KANJIDIC2 file which may not always be the most important one. For full readings/meanings use the function [lookup](#) or consult a dictionary.

## Source

Most of the data is directly from the KANJIDIC2 file. [https://www.edrdg.org/wiki/index.php/KANJIDIC\\_Project](https://www.edrdg.org/wiki/index.php/KANJIDIC_Project)

Variables **jlpt**, **frank**, **idc**, **components** were taken from the Kanjium data base <https://github.com/mifunetoshiro/kanjium>

Variable **components** is originally from RADKFILE/KRADFILE. <https://www.edrdg.org/>

The use of this data is covered in each case by a Creative Commons BY-SA 4.0 License. See the package's LICENSE file for details and copyright holders.

Variable "class" is derived from "grade".

Variable "kanken" was compiled based on the Wikipedia description of the test levels (as of September 2022).

---

<b>kanjidist</b>	<i>Compute distance between two kanjivec objects based on hierarchical optimal transport</i>
------------------	--

---

## Description

The kanji distance is based on matching hierarchical component structures in a nesting-free way across all levels. The cost for matching individual components is a cost for registering the components (i.e. aligning their position, scale and aspect ratio) plus the (relative unbalanced) Wasserstein distance between the registered components.

## Usage

```
kanjidist(
  k1,
  k2,
  compo_seg_depth1 = 3,
  compo_seg_depth2 = 3,
  p = 1,
  C = 0.2,
  approx = c("grid", "pc", "pcweighted"),
  type = c("rtt", "unbalanced", "balanced"),
  size = 48,
  lwd = 2.5,
  density = 30,
  verbose = FALSE,
  minor_warnings = TRUE
)
```

## Arguments

<b>k1, k2</b>	two objects of type <code>kanjivec</code> .
<b>compo_seg_depth1, compo_seg_depth2</b>	two integers $\geq 1$ . Specifies for each kanji the deepest level included for component matching. If 1, only the kanji itself is used.
<b>p</b>	the order of the Wasserstein distance used for matching components. All distances and the penalty (if any) are taken to the $p$ -th power (which is compensated by taking the $p$ -th root after summation).
<b>C</b>	the penalty for extra mass if <code>type</code> is <code>"rtt"</code> or <code>"unbalanced"</code> , i.e. we add $C^p$ per unit of extra mass (before applying the $p$ -th root).
<b>approx</b>	what kind of approximation is used for matching components. If this is <code>"grid"</code> , a bitmap (raster image) is used, otherwise lines are approximated by more freely spaced points. For <code>"pc"</code> (point cloud) each point has the same weight and points are placed in a (more or less) equidistant way. For <code>"pcweighted"</code> points are further apart along straight lines and around the center of the Bezier curves that describe the strokes. The weights of the

points are then (more or less) proportional to the amount of ink (stroke length) they represent.

<code>type</code>	the type of Wasserstein distance used for matching components based on the grid or point cloud approximation chosen. "unbalanced" means the weights (pixel values if <code>approx = "grid"</code> ) are interpreted as mass. The total masses in two components be very different. Extra mass can be disposed of at cost $C^p$ per unit. "rtt" is computationally the same, but the final distance is divided by the maximum of the total ink (sum of weights) in each component to the $1/p$ . "balanced" means the weights are normalized so that both images have the same total mass 1. Everything has to be transported, i.e. disposal of mass is not allowed.
<code>size</code>	side length of the bitmaps used for matching components (if <code>approx = "grid"</code> ).
<code>lwd</code>	linewidth for drawing the components in these bitmaps (if <code>approx = "grid"</code> ).
<code>density</code>	approximate number of discretization points per unit line length (if <code>approx != "grid"</code> )
<code>verbose</code>	logical. Whether to print detailed information on the cost for all pairs of components and the final matching.
<code>minor_warnings</code>	logical. Should <code>minor_warnings</code> be given. If <code>FALSE</code> , the warnings about substantial distances between bitmaps/pointclouds standing for the same component and the use of a workaround due to missing strokes in component decompositions are suppressed. While these warnings indicate to some extent that things are not going exactly as planned, they are usually not of interest if a larger number of kanji distances is computed and obscure the visibility of more important warnings (if any).

## Details

For the precise definition and details see the reference below. Parameter `C` corresponds to  $b/2^{1/p}$  in the paper.

## Value

The kanji distance, a non-negative number.

## Warning

### [Experimental]

The interface and details of this function will change in the future. Currently only a minimal set of parameters can be passed. The other parameters are fixed exactly as in the "prototype distance" (4.1) of the reference below for better or worse.

There is a certain tendency that exact matches of components are rather strongly favored (if the KanjiVG elements agree this can overrule the unbalanced Wasserstein distance) and the penalties for translation/scaling/distortion of components are somewhat mild.

The computation time is rather high (depending on the settings and kanji up to several seconds per kanji pair). This can be alleviated somewhat by keeping the `compo_seg_depth` parameters at 3 or lower and setting `size = 32` (which goes well with `lwd=1.8`).

Future versions will use a much faster line base optimal transport algorithm and further speed-ups.

## References

Dominic Schuhmacher (2023).  
 Distance maps between Japanese kanji characters based on hierarchical optimal transport.  
 ArXiv, [doi:10.48550/arXiv.2304.02493](https://doi.org/10.48550/arXiv.2304.02493)

## See Also

[kanjidistmat](#), [kmatdist](#)

## Examples

```
if (requireNamespace("ROI.plugin.glpk")) {
  kanjidist(fivebetas[[4]], fivebetas[[5]])
  kanjidist(fivebetas[[4]], fivebetas[[5]], verbose=TRUE)
  # faster and similar:
  kanjidist(fivebetas[[4]], fivebetas[[5]], compo_seg_depth1=2, compo_seg_depth2=2,
            size=32, lwd=1.8, verbose=TRUE)
  # slower and similar:
  kanjidist(fivebetas[[4]], fivebetas[[5]], size=64, lwd=3.2, verbose=TRUE)
}
```

---

kanjidistmat	<i>Compute distance matrix based on hierarchical optimal transport for lists of <code>kanjivec</code> objects</i>
--------------	---

---

## Description

Individual distances are based on [kanjidist](#).

## Usage

```
kanjidistmat(
  klist,
  klist2 = NULL,
  compo_seg_depth = 3,
  p = 1,
  C = 0.2,
  approx = c("grid", "pc", "pcweighted"),
  type = c("rtt", "unbalanced", "balanced"),
  size = 48,
  lwd = 2.5,
  density = 30,
  verbose = FALSE,
  minor_warnings = FALSE
)
```

**Arguments**

- `klist` a list of [kanjimat](#) objects.
- `klist2` an optional second list of [kanjimat](#) objects.
- `compo_seg_depth` integer  $\geq 1$ . Specifies for all kanji the deepest level included for component matching. If 1, only the kanji itself is used.
- `p, C, type, approx, size, lwd, density, verbose, minor_warnings` the same as for the function [kanjidist](#), with the sole difference that `minor_warnings` defaults to `FALSE` here.

**Value**

A matrix of dimension `length(klist) x length(klist2)` having as its  $(i, j)$ -th entry the distance between `klist[[i]]` and `klist2[[j]]`. If `klist2` is not provided it is assumed to be equal to `klist`, but computation is more efficient as only the upper triangular part is computed and then symmetrized with diagonal zero.

**Warning**

**[Experimental]**  
The same precautions apply as for [kanjidist](#).

**See Also**

[kanjidist](#), [kmatdistmat](#)

**Examples**

```
kanjidistmat(fivebetas)
```

---

`kanjimat`

*Create kanjimat objects*

---

**Description**

Create a (list of) `kanjimat` object(s), i.e. bitmap representations of a kanji using a certain font-family and other typographical parameters.

**Usage**

```
kanjimat(  
  kanji,  
  family = NULL,  
  size = NULL,  
  margin = 0,
```

```

    antialias = TRUE,
    save = FALSE,
    overwrite = FALSE,
    simplify = TRUE,
    ...
)

```

### Arguments

<code>kanji</code>	a (vector of) character string(s) containing kanji.
<code>family</code>	the font-family to be used. For details see vignette.
<code>size</code>	the sidelength of the (square) bitmap
<code>margin</code>	numeric. Extra margin around the character. Defaults to 0 which leaves a relatively slim margin. Positive values increase this margin, negative values decrease it (which usually cuts off part of the kanji).
<code>antialias</code>	logical. Shall antialiasing be performed?
<code>save</code>	logical or character. If <code>FALSE</code> return the (list of) <code>kanjimat</code> object(s). Otherwise save the result as an rds file in the working directory (as <code>kmatsave.rds</code> ) or under the file path provided.
<code>overwrite</code>	logical. If <code>FALSE</code> return an error (before any computations are done) if the designated file path already exists. Otherwise an existing file is overwritten.
<code>simplify</code>	logical. Shall a single <code>kanjimat</code> object be returned (instead a list of one) if <code>kanji</code> is a single kanji?
<code>...</code>	further arguments passed to <a href="#">png</a> . This is for extensibility. The only argument that may currently be used is <code>type</code> . Trying to change sizes, units, colors or fonts by this argument results in an error or an undesirable output.

### Value

A list of objects of class `kanjimat` or, if only one kanji was specified and `simplify` is `TRUE`, a single objects of class `kanjimat`. If `save = TRUE`, the same is (saved and) still returned invisibly.

### Warning

If no font family is provided, the default **Chinese** font WenQuanYi Micro Hei that comes with the package showtext is used. This means that the characters will typically be recognizable, but quite often look odd as Japanese characters. We strongly advised that a Japanese font is used as detailed above.

### Examples

```
res <- kanjimat(kanji=" ", size = 128)
```



---

<code>kanjivec</code>	<i>Create kanjivec objects from kanjivg data</i>
-----------------------	--

---

## Description

Create a (list of) kanjivec object(s). Each object is a representation of the kanji as a tree of strokes based on .svg files from the KanjiVG database containing further, derived information.

## Usage

```
kanjivec(
  kanji,
  database = NULL,
  flatten = "intelligent",
  bezier_discr = c("svgparser", "eqtimed", "eqspaced"),
  save = FALSE,
  overwrite = FALSE,
  simplify = TRUE
)
```

## Arguments

<code>kanji</code>	a (vector of) character string(s) of one or several kanji.
<code>database</code>	the path to a local copy of (a subset of) the KanjiVG database. It is expected that the svg files reside at this exact location (not in a sub-directory). If <code>NULL</code> , an attempt is made to read the svg file(s) from the KanjiVG GitHub repository (after prompting for confirmation, which can be switched off via the <a href="#">option ask_github</a> ).
<code>flatten</code>	logical. Should nodes that are only-children be fused with their parents? Alternatively one of the strings "intelligent", "inner" or "leaves". Although the first is the default it is experimental and the precise meaning will change in the future; see details.
<code>bezier_discr</code>	character. How to discretize the Bézier curves describing the strokes. If "svgparser" (the only option available prior to kanjivat 0.12.0), code from the non-CRAN package <code>svgparser</code> is used for discretizing at equal time steps. The new choices "eqtimed" and "eqspaced" discretize into fewer points (and allow for more customization underneath). The former creates discretization points at equal time steps, the latter at equal distance steps (to a good approximation).
<code>save</code>	logical or character. If <code>FALSE</code> return the (list of) kanjivec object(s). Otherwise save the result as an rds file in the working directory (as <code>kvec-save.rds</code> ) or under the file path provided.
<code>overwrite</code>	logical. If <code>FALSE</code> return an error (before any computations are done) if the designated file path already exists. Otherwise an existing file is overwritten.

**simplify** logical. Shall a single kanjivec object be returned (instead a list of one) if `kanji` is a single kanji?

## Details

A kanjivec object contains detailed information on the strokes of which an individual kanji is composed including their order, a segmentation into reasonable components ("radicals" in a more general sense of the word), classification of individual strokes, and both vector data and interpolated points to recreate the actual stroke in a Kyoukashou style font. For more information on the original data see <http://kanjivg.tagaini.net/>. That data is licenced under Creative Commons BY-SA 3.0 (see licence file of this package).

The original .svg files sometimes contain additional `<g>` elements that provide information about the current group of strokes rather than establishing a new subgroup of its own. This happens typically for information that establishes coherence with another part of the tree (by noting that the current subgroup is also part 2 of something else), but also for variant information. With the option `flatten = TRUE` the extra hierarchy level in the tree is avoided, while the original information in the KanjiVG file is kept. This is achieved by fusing only-children to their parents, giving the new node the name of the child and all its attributes, but prefixing `p.` to the attribute names of the parent (the parents' "names" attribute is discarded, but can be reconstructed from the parents' id). Removal of several hierarchies in sequence can lead to attribute names with multiple `p.` in front. Fusing to parents is suppressed if the parent is the root of the hierarchy (typically for one-stroke kanji), as this could lead to confusing results.

The options `flatten = "inner"` and `flatten = "leaves"` implement the above behavior only for the corresponding type of node (inner nodes or leaves). The option `flatten = "intelligent"` tries to find out in more sophisticated ways which flattening is desirable and which is not (it will flatten rather conservatively). Currently nodes without an element attribute that have only one child are flattened away (one example where this is reasonable is in kanji `kbase[187, ]`), as are nodes with an element attribute and only one child if this child is also an inner node and has the same element and part attribute as the parent, but both have no number (this would be problematic for any component-building code in the particular case of kanji `kbase[1111, ]`).

A `kanjivec` object has components

**char** the kanji (a single character)

**hex** its Unicode codepoint (integer of class `hexmode`)

**padhex** the Unicode codepoint padded with zeros to five digits (mode character)

**family** the font on which the data is based. Currently only "schoolbook" (to be extended with "kaisho" at some point)

**nstrokes** the number of strokes in the kanji

**ncompos** a vector of the number of components at each depth of the tree

**nveins** the number of veins in the component structure

**strokedend** the decomposition tree of the kanji as an object of class `dendrogram`

**components** the component structure by segmentation depth (components can overlap) in terms of KanjiVG elements and their depth-first tree coordinates

**veins** the veins in the component structure. Each vein is represented as a two-column matrix that lists in its rows the indices of **components** (starting at the root, which in the component indexing is `c(1,1)`)

**stroketree** the decomposition tree of the kanji, a list containing the full information of the the KanjiVG file (except some top level attributes)

**stroketree** is a close representation of the KanjiVG svg file as list object with some serious nesting of sublists. The XML attributes become attributes of the list and its elements. The user will usually not have to look at or manipulate **stroketree** directly, but **strokedend** and **compents** are derived from it and other functions may process it further.

The main differences to the svg file are

1. the actual strokes are not only given as d-attributes describing Bézier curves, but but also as two-column matrices describing discretizations of these curves. These matrices are the actual contents of the innermost lists in **stroketree**, but are more conveniently accessed via the function `get_strokes`. Starting with version 0.13.0, there is also an additional attribute "beziermat", which describes the Bézier curves for the stroke in a  $2 \times (1+3n)$  matrix format. The first column is the start point, then each triplet of columns stands for control point 1, control point 2 and end point (=start point of the next Bézier curve if any).
2. The positions of the stroke numbers (for plotting) are saved as an attribute `strokenum_coords` to the entire stroke tree rather than a separate element.

**strokedend** is more easy to examine and work with due to various convenience functions for dendrograms in the packages **stats** and **dendextend**, including `str` and `plot.dendrogram`. The function `plot.kanjivec` with option `type = "dend"` is a wrapper for `plot.dendrogram` with reasonable presets for various options.

The label-attributes of the nodes of **strokedend** are taken from the element (for inner nodes) and `type` (for leaves) attributes of the .svg files. They consist of UTF-8 characters representing kanji parts and a combination of UTF-8 characters for representing strokes and may not represent well in all CJK fonts (see details of `plot.kanjivec`). If element and `type` are missing in the .svg file, the label assigned is the second part of the id-attribute, e.g. `g5` or `s9`.

The **components** at a given level can be plotted, see `plot.kanjivec` with `type = "kanji"`. Both **components** and **veins** serve mainly for the computation of **kanji distances**.

## Value

A list of objects of class **kanjivec** or, if only one kanji was specified and `simplify` is `TRUE`, a single objects of class **kanjivec**. If `save = TRUE`, the same is (saved and) still returned invisibly.

## See Also

`plot.kanjivec`, `str.kanjivec`

## Examples

```
if (interactive()) {
  # Try to load the svg file for the kanji from GitHub.
  res <- kanjivec(" ", database=NULL)
  str(res)
}

fivebetas # sample kanjivec data
str(fivebetas[[1]])
```

---

<code>kmatdist</code>	<i>Compute the unbalanced or balanced Wasserstein distance between two kanjimat objects</i>
-----------------------	---

---

## Description

This gives the dissimilarity of pixel-images of the kanji based on how far mass (or "ink") has to be transported to transform one image into the other.

## Usage

```
kmatdist(
  k1,
  k2,
  p = 1,
  C = 0.2,
  type = c("unbalanced", "balanced"),
  output = c("dist", "all")
)
```

## Arguments

<code>k1, k2</code>	two objects of type <code>kanjimat</code> .
<code>p</code>	the order of the Wasserstein distance. All distances and a potential penalty are taken to the <code>p</code> -th power (which is compensated by taking the <code>p</code> -th root after summation).
<code>C</code>	the penalty for extra mass if <code>type="unbalanced"</code> , i.e. we add $C^p$ per unit of extra mass (before applying the <code>p</code> -th root).
<code>type</code>	the type of Wasserstein metric. <code>"unbalanced"</code> means the pixel values in the two images are interpreted as mass. The total masses can be very different. Extra mass can be disposed of at cost $C^p$ per unit. <code>"balanced"</code> means the pixel values are normalized so that both images have the same total mass 1. Everything has to be transported, i.e. disposal of mass is not allowed.
<code>output</code>	the requested output. See return value below.

**Value**

If `output = "dist"`, a single non-negative number: the unbalanced or balanced Wasserstein distance between the kanji. If `output = "all"` a list with detailed information on the transport plan and the disposal of pixel mass. See [unbalanced](#) for details.

**See Also**

[kmatdistmat](#), [kanjidist](#)

**Examples**

```
res <- kmatdist(fivetrees1[[1]], fivetrees1[[5]], p=1, C=0.1, output="all")
plot(res, what="plan", angle=20, lwd=1.5)
plot(res, what="trans")
plot(res, what="extra")
plot(res, what="inplace")
```

---

**kmatdistmat**
*Compute distance matrix for lists of kanjimat objects*


---

**Description**

Apply [kmatdist](#) to every pair of [kanjimat](#) objects to compute the unbalanced or balanced Wasserstein distance.

**Usage**

```
kmatdistmat(
  klist,
  klist2 = NULL,
  p = 1,
  C = 0.2,
  type = c("unbalanced", "balanced")
)
```

**Arguments**

`klist` a list of [kanjimat](#) objects.  
`klist2` an optional second list of [kanjimat](#) objects.  
`p, C, type` the same as for the function [kmatdist](#).

**Value**

A matrix of dimension `length(klist) x length(klist2)` having as its  $(i, j)$ -th entry the distance between `klist[[i]]` and `klist2[[j]]`. If `klist2` is not provided it is assumed to be equal to `klist`, but the computation is more efficient as only the upper triangular part is computed and then symmetrized with diagonal zero.

**See Also**

[kmatdist](#), [kanjidistmat](#)

**Examples**

```
kmatdistmat(fivetrees1)
kmatdistmat(fivetrees1, fivetrees1) # same result but slower
kmatdistmat(fivetrees1, fivetrees2) # note the smaller values on the diagonal
```

---

`kreadmean`

*Kanji readings and meanings*

---

**Description**

Data set of all kanji readings and meanings from the KANJIDIC2 dataset in an R list format. For convenient access to this data use function [lookup](#).

**Usage**

```
kreadmean
```

**Format**

An object of class `list` of length 13108.

**Source**

KANJIDIC2 file by Jim Breen and The Electronic Dictionary Research and Development Group (EDRDG)

[https://www.edrdg.org/wiki/index.php/KANJIDIC\\_Project](https://www.edrdg.org/wiki/index.php/KANJIDIC_Project)

The use of this data is covered by the Creative Commons BY-SA 4.0 License.

---

`lookup`

*Look up kanji*

---

**Description**

Return readings and meanings or information from `kbase` or `kmorph`.

**Usage**

```
lookup(kanji, what = c("readmean", "basic", "morphologic"))
```

**Arguments**

**kanji** a (vector of) character strings containing kanji.  
**what** the sort of information to display.

**Details**

This is a very basic interface for a quick lookup information based on exact knowledge of the kanji (provided by a Japanese input method or its UTF-8 code). Most of the information is based on the KANJIDIC2 file by EDRDG (see thank you page) Please use one of the many excellent online kanji dictionaries (see e.g.) more sophisticated lookup methods and more detailed results.

**Value**

If **what** is "readmean" the information is output with cat and there is no return value (invisible NULL) In the other cases the appropriate subsets of the tables kbase and kmorph are returned

**Author(s)**

Dominic Schuhmacher <schuhmacher@math.uni-goettingen.de>

**Examples**

```
lookup(c(" ", " ", " "))
lookup(" ") # same
```

---

<code>options</code>	<i>Kanjistat Options</i>
----------------------	--------------------------

---

**Description**

Set or examine global kanjistat options.

**Usage**

```
kanjistat_options(...)  

get_kanjistat_option(x)
```

**Arguments**

**...** any number of options specified as **name = value**  
**x** name of an option given as character string.

**Value**

`kanjimat_options` returns the list of all set options if there is no function argument. Otherwise it returns list of *all* old options. `get_kanjimat_option` returns the current value set for option `x` or NULL if the option is not set.

---

<code>plot.kanjimat</code>	<i>Plot kanjimat object</i>
----------------------------	-----------------------------

---

**Description**

Plot kanjimat object

**Usage**

```
## S3 method for class 'kanjimat'
plot(
  x,
  mode = c("dark", "light"),
  col = gray(seq(0, 1, length.out = 256)),
  ...
)
```

**Arguments**

<code>x</code>	object of class <code>kanjimat</code> .
<code>mode</code>	character string. If "dark" the original grayscale values are used, if "light" they are inverted. With the default grayscale color scheme the kanji is plotted white-on-black for "dark" and black-on-white for "light".
<code>col</code>	a vector of colors. Typically 256 values are enough to keep the full information of an (antialiased) <code>kanjimat</code> object.
<code>...</code>	further parameters passed to <code>image</code> .

**Value**

No return value, called for side effects.



---

plot.kanjivec	<i>Plot kanjivec objects</i>
---------------	------------------------------

---

## Description

Plot kanjivec objects

## Usage

```
## S3 method for class 'kanjivec'
plot(
  x,
  type = c("kanji", "dend"),
  seg_depth = 0,
  palette = "Dark 3",
  pal.extra = 0,
  numbers = FALSE,
  offset = c(0.025, 0),
  family = NULL,
  lwd = 8,
  ...
)
```

## Arguments

<b>x</b>	an object of class <code>kanjivec</code>
<b>type</b>	either "kanji" or "dend". Whether to plot the actual kanji, coloring strokes according to levels of segmentation, or to plot a representation of the tree structure underlying this segmentation. Among the following named parameters, only <b>family</b> is for use with <b>type = "dend"</b> ; all others are for <b>type = "kanji"</b> .
<b>seg_depth</b>	an integer. How many steps down the segmentation hierarchy we use different colors for different groups. If zero (the default), only one color is used that can be specified with <b>col</b> passed via <b>...</b> as usual
<b>palette</b>	a valid name of a hcl palette (one of <code>hcl.pals()</code> ). Used for coloring the components if <b>seg_depth</b> is $> 0$ .
<b>pal.extra</b>	an integer. How many extra colors are picked in the specified palette. If this is 0 (the default), <b>palette</b> is used with as many colors as we have components. Since many hcl palettes run from dark to light colors, the last (few) components may be too light. Increasing <b>pal.extra</b> then makes the component colors somewhat more similar, but the last component darker.
<b>numbers</b>	logical. Shall the stroke numbers be displayed.
<b>offset</b>	the (x,y)-offset for the numbers relative to the positions from <code>kanjivg</code> saved in the <code>kanjivec</code> object. Either a vector of length 2 specifying some

	fixed offset for all numbers or a matrix of dimension <code>kanjivec\$nstrokes</code> times 2.
<code>family</code>	the font-family for labeling the nodes if <code>type = "dend"</code> . See details.
<code>lwd</code>	the usual line width graphics parameter.
<code>...</code>	further parameters passed to <code>lines</code> if <code>type = "kanji"</code> and to <code>plot.dendrogram</code> if <code>type = "dend"</code> .

### Details

Setting up nice labels for the nodes if `type = "dend"` is not easy. For many font families it appears that some "kanji components" cannot be displayed in plots even with the help of package `showtext` and if the font contains glyphs for the corresponding codepoints that display correctly in text documents. This concerns in increasing severity of the problem Unicode blocks 2F00–2FDF (Kangxi Radicals), 2E80–2EFF (CJK Radicals Supplement) and 31C0–31EF (CJK Strokes). For the strokes it seems nearly impossible which is why leaves are simply annotated with the number of the strokes.

For the other it is up to the user to find a suitable font and pass it via the argument `font.family`. The default `family = NULL` first tries to use `default_font` if this option has been set (via `kanjistat_options`) and otherwise uses `wqy-microhei`, the Chinese default font that comes with package `showtext` and cannot display any radicals from the supplement.

On a Mac the experience is that "hiragino\_sans" works well. In addition there is the issue of font size which is currently not judiciously set and may be too large for some (especially on-screen) devices. The parameter `cex` (via `...`) fixes this.

### Value

No return value, called for side effects.

### Examples

```
kanji <- fivebetas[[2]]
plot(kanji, type = "kanji", seg_depth = 2)
plot(kanji, type = "dend")
# gives a warning if get_kanjistat_option("default_font") is NULL
```

---

plotkanji

*Plot kanji*

---

### Description

Write kanji to a graphics device.

**Usage**

```
plotkanji(
  kanji,
  device = "default",
  family = NULL,
  factor = 10,
  width = NULL,
  height = NULL,
  ...
)
```

**Arguments**

<code>kanji</code>	a vector of class character specifying one or several kanji to be plotted.
<code>device</code>	the type of graphics device where the kanji is plotted. Defaults to the user's default type according to <code>getOption("device")</code> .
<code>family</code>	the font family or families used for writing the kanji. Make sure to add the font(s) first by using <code>font_add</code> ; see details. If <code>family</code> is a vector of several font families they are matched to the characters in <code>kanji</code> (and possibly recycled).
<code>factor</code>	a magnification factor applied to the font size (typically 12 points).
<code>width, height</code>	the dimensions of the device.
<code>...</code>	further parameters passed to the function opening the device (such as a file name for devices that create a file).

**Details**

This function writes one or several kanji to a graphics device in an arbitrary font that has been registered, i.e., added to the database in package `sysfonts`. For the latter say `font_add` or `font_families` to verify what fonts are available.

For further information see *Working with Japanese fonts* in `vignette("kanjistat", package = "kanjistat")`. `plotkanji` uses the package `showtext` to write the kanji in a large font at the center of a new device of the specified type. specify `device = "current"` to write the kanji to the current device. It is now recommended to simply use `graphics::text` in combination with `showtext::showtext_auto` instead.

**Value**

No return value, called for side effects.

**Warning**

If no font family is provided, the default **Chinese** font WenQuanYi Micro Hei that comes with the package `showtext` is used. This means that the characters will typically be recognizable, but quite often look odd as Japanese characters. We strongly advised that a Japanese font is used as detailed above.

**Examples**

```
plotkanji(" ")
plotkanji(" ")
```

---

pooled\_similarity      *Precomputed kanji distances*

---

**Description**

Precomputed kanji distances

**Usage**

```
pooled_similarity
```

**Format**

A tibble containing kanji similarity judgments by 3 "native or native-like" speakers of Japanese. For each row, the pivot kanji was compared to a list of potential distractors. From the distractors, the subjects selected one character which they found particularly easy to confuse with the pivot. For the exact methodology, see the original study referenced below.

**Source**

Datasets from <https://lars.yencken.org/datasets>, made available under the Creative Commons Attribution 3.0 Unported licence.

Collected as part of *Yencken, Lars (2010) Orthographic support for passing the reading hurdle in Japanese. PhD Thesis, University of Melbourne, Melbourne, Australia.*

**References**

Yencken, Lars, & Baldwin, Timothy (2008). Measuring and predicting orthographic associations: Modelling the similarity of Japanese kanji. In: *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pp. 1041-1048.

**Examples**

```
# Get kanji characters that were found to be easily confused with .
pooled_similarity[pooled_similarity$selected == " ", ]$pivot
```

---

<code>print.kanjivec</code>	<i>Print basic information about a kanjivec object</i>
-----------------------------	--

---

### Description

Print basic information about a kanjivec object

### Usage

```
## S3 method for class 'kanjivec'
print(x, dend = FALSE, ...)
```

### Arguments

<code>x</code>	an object of class <code>kanjivec</code> .
<code>dend</code>	whether to print the structure of the <code>strokedend</code> component.
<code>...</code>	further parameters passed to <code>print.default</code> .

### Value

No return value, called for side effects.

---

<code>read_kanjidic2</code>	<i>Read a KANJIDIC2 file</i>
-----------------------------	------------------------------

---

### Description

Perform basic validity checks and transform data to a standardized list or keep as an object of class `xml_document` (package `xml2`).

### Usage

```
read_kanjidic2(fpath = NULL, output = c("list", "xml"))
```

### Arguments

<code>fpath</code>	the path to a local KANJIDIC2 file. If <code>NULL</code> (the default) the most recent KANJIDIC2 file is downloaded from <a href="https://www.edrdg.org/kanjidic/kanjidic2.xml.gz">https://www.edrdg.org/kanjidic/kanjidic2.xml.gz</a> after asking for confirmation.
<code>output</code>	one of "list" or "xml". The desired type of output.

## Details

KANJIDIC2 contains detailed information on all of the 13108 kanji in three main Japanese standards (JIS X 0208, 0212 and 0213). The KANJIDIC files have been compiled and maintained by Jim Breen since 1991, with the help of various other people. The copyright is now held by the Electronic Dictionary Research and Development Group (EDRDG). The files are made available under the Creative Commons BY-SA 4.0 license. See [https://www.edrdg.org/wiki/index.php/KANJIDIC\\_Project](https://www.edrdg.org/wiki/index.php/KANJIDIC_Project) for details on the contents of the files and their license.

If `output = "xml"`, some minimal checks are performed (high level structure and total number of kanji).

If `output = "list"`, additional validity checks of the lower level structure are performed. Most are in accordance with the file's Document Type Definition (DTD). Some additional check concern some common patterns that are true about the current KANJIDIC2 file (as of December 2023) and seem unlikely to change in the near future. This includes that there is always at most one `rmgroup` entry in `reading_meaning`. Informative warnings are provided if any of these additional checks fail.

## Value

If `output = "xml"`, the exact XML document obtained from `xml2::read_xml`. If `output = "list"`, a list of lists (the individual kanji), each with the following seven components.

- **literal**: a single UTF-8 character representing the kanji.
- **codepoint**: a named character vector giving the available codepoints in the unicode and jis standards.
- **radical**: a named numeric vector giving the radical number(s), in the range 1 to 214. The number named `classical` is as recorded in the *KangXi Zidian* (1716); if there is a number named `nelson_c`, the kanji was reclassified in Nelson's *Modern Reader's Japanese-English Character Dictionary* (1962/74).
- **misc**: a list with six components
  - **grade**: the kanji grade level. 1 through 6 indicates a kyouiku kanji and the grade in which the kanji is taught in Japanese primary school. 8 indicates one of the remaining jouyou kanji learned in junior high school, and 9 or 10 are jinmeiyou kanji. The remaining (hyougai) kanji have NA as their entry.
  - **stroke\_count**: The stroke count of the kanji, including the radical. If more than one, the first is considered the accepted count, while subsequent ones are common miscounts.
  - **variant**: a named character vector giving either a cross-reference code to another kanji, usually regarded as a variant, or an alternative indexing code for the current kanji. The type of variant is given in the name.
  - **freq**: the frequency rank (1 = most frequent) based on newspaper data. NA if not among the 2500 most frequent.
  - **rad\_name**: a character vector. For a kanji that is a radical itself, the name(s) of the radical (if there are any), otherwise of length 0.
  - **jlpt**: The Japanese Language Proficiency Test level according to the old four-level system that was in place before 2010. A value from 4 (most elementary) to 1 (most advanced).

- **dic\_number**: a named character vector (possibly of length 0) giving the index numbers (for some kanji with letters attached) of the kanji in various dictionaries, textbooks and flashcard collections (specified by the name). For Morohashi's *Dai Kan-Wa Jiten*, the volume and page number is also provided in the format `moro.VOL.PAGE`.
- **query\_code**: a named character vector giving the codes of the kanji in various query systems (specified by the name). For Halpern's SKIP code, possible misclassifications (if any) of the kanji are also noted in the format `mis.skip.TYPE`, where TYPE indicates the type of misclassification.
- **reading\_meaning**: a (possibly empty) list containing zero or more **rmgroup** components creating groups of readings and meanings (in practice there is never more than one **rmgroup** currently) as well as a component **nanori** giving a character vector (possibly of length 0) of readings only associated with names. Each **rmgroup** is a list with entries:
  - **reading**: a (possibly empty) list of entries named from among `pinyin`, `korean_r`, `korean_h`, `vietnam`, `ja_on` and `ja_kun`, each containing a character vector of the corresponding readings
  - **meaning**: a (possibly empty) list of entries named with two-letter (ISO 639-1) language codes, each containing a character vector of the corresponding meanings.

## See Also

[kanjidata](#), [kreadmean](#)

## Examples

```
if (interactive()) {
  read_kanjidic2("kanjidic2.xml")
}
```

---

samplekan

*Sample kanji from a set*

---

## Description

Sample kanji from a set

## Usage

```
samplekan(
  set = c("kyouiku", "jouyou", "jinmeiyou", "kanjidic"),
  size = 1,
  replace = FALSE,
  prob = NULL
)
```

**Arguments**

<b>set</b>	a character string specifying the set of kanjis to sample from.
<b>size</b>	a positive number, the number of samples.
<b>replace</b>	logical. Sample with replacement?
<b>prob</b>	currently without effect.

**Value**

a vector of length **size** containing the individual characters

**Examples**

```
(sam <- samplekan(size = 10))
lookup(sam)
```

---

<b>sedist</b>	<i>Compute the stroke edit distances between two sets of kanji</i>
---------------	--

---

**Description**

Variants of the stroke edit distance proposed by Yencken (2010). Each kanji is encoded as sequence of stroke types according to its stroke order, using the type attribute from the kanjiVG data. Then the edit distance (a.k.a. \ Levenshtein distance) between sequences is computed and divided by the maximum of the number of strokes

**Usage**

```
sedist(k1, k2, type = c("full", "before_slash", "first"))
```

**Arguments**

<b>k1, k2</b>	atomic vectors or lists of kanji in any format that can be treated by <a href="#">convert_kanji()</a>
<b>type</b>	the type of stroke edit distance to compute. See details.

**Details**

The kanjiVG type attribute is a single string composed of a CJK strokes Unicode character, an optional latin letter providing further information and possibly a variant (another CJK strokes character with optional letter) separated by "/". If **type** is "full" a match is only counted if two strings are exactly the same, "before\_slash" ignores any slashes and what comes after them, "first" only considers the first character of each string (so the first CJK stroke character) when counting matches.

The stroke edit distance used by Yencken (2010) is obtained by setting **type** = "all" (the default), except that the underlying kanjiVG data has significantly changed since then. Comparing with the values in [dstrokedist](#) we get an agreement of 96.3 percent, whereas the other distances disagree by a small amount (usually 1-2 edit operations).



**Value**

A `length(k1) x length(k2)` matrix of stroke edit distances.

**Warning**

Requires `kanjistat.data` package.

**References**

Yencken, Lars (2010). Orthographic support for passing the reading hurdle in Japanese. PhD Thesis, University of Melbourne, Australia

**Examples**

```
ind1 <- 384
k1 <- convert_kanji(ind1, "character")
ind2 <- which(dstrokededit[ind1,] > 0)
# dstrokededit contains only the "closest" kanji
k2 <- convert_kanji(ind2, "character")
row_a <- dstrokededit[ind1, ind2]
if (requireNamespace("kanjistat.data", quietly = TRUE)) {
  row_b <- sedist(k1, k2)
  mat <- rbind(row_a, row_b)
  rownames(mat) = c(k1, k1)
  colnames(mat) = k2
  mat
}
```

---

`str.kanjivec`

*Compactly display the structure of a kanjivec object*

---

**Description**

Compactly display the structure of a `kanjivec` object

**Usage**

```
## S3 method for class 'kanjivec'
str(object, ...)
```

**Arguments**

`object` an object of class `kanjivec`.  
`...` further parameters passed to `str` for all but the `stroketree` component of `object`.

**Value**

No return value, called for side effects.

# Index

- \* datasets
  - distdata, 6
  - fivebetas, 7
  - fivetrees, 8
  - kanjidata, 10
  - kreadmean, 22
  - pooled\_similarity, 28
- cjk\_escape, 2
- codepoint, 3
- codepointToKanji (*codepoint*), 3
- compare\_neighborhoods, 4
- convert\_kanji, 5
- convert\_kanji(), 32
- dendextend, 19
- distdata, 6
- dstrokedit, 32
- dstrokedit (*distdata*), 6
- dyehli (*distdata*), 6
- fivebetas, 7
- fivetrees, 8
- fivetrees1 (*fivetrees*), 8
- fivetrees2 (*fivetrees*), 8
- fivetrees3 (*fivetrees*), 8
- font\_add, 27
- font\_families, 27
- get\_kanjistat\_option (*options*), 23
- get\_strokes, 9, 10, 19
- get\_strokes\_compo, 9, 9
- image, 24
- kanji distances, 19
- kanjidata, 10, 31
- kanjidist, 12, 14, 15, 21
- kanjidist(), 4
- kanjidistmat, 14, 14, 22
- kanjimat, 8, 15, 15, 21
- kanjimat\_options, 26
- kanjimat\_options (*options*), 23
- kanjiToCodepoint (*codepoint*), 3
- kanjivec, 5, 7, 17, 25
- kbase, 5, 6
- kbase (*kanjidata*), 10
- kmatdist, 8, 14, 20, 21, 22
- kmatdistmat, 15, 21, 21
- kmorph (*kanjidata*), 10
- kreadmean, 22, 31
- lookup, 11, 22, 22
- option, 17
- options, 23
- plot.dendrogram, 19
- plot.kanjimat, 24
- plot.kanjivec, 19, 25
- plotkanji, 26
- png, 16
- pooled\_similarity, 28
- print.kanjivec, 29
- read\_kanjidic2, 29
- samplekan, 31
- sedist, 32
- str, 19
- str.kanjivec, 19, 33
- unbalanced, 21
- xml2::read\_xml, 30
- xml\_document, 29