

functions.R

```
1 ##' Inserts escape characters to protect preserved symbols in LaTeX
2 ##' markup
3 ##'
4 ##' Found fixes for all of the LaTeX symbols except backslash
5 ##' @param x a string
6 ##' @return corrected character vector for LaTeX
7 ##' @author Paul Johnson
8 ##' @examples
9 ##' x <- c("_asdf&_&$", "asd adf asd_", "~ % & $asdf_")
10 ##' escape(x)
11 escape <- function (x){
12   x <- gsub("\\\\", "SANITIZE.BACKSLASH", x)
13   for(i in c("&", "$", "_", "^", "%", "#", "{", "{")) {
14     x <- gsub(paste0("\\", i), paste0("\\\\",i), x)
15   }
16   x <- gsub("\\~", "$\\sim$", x)
17   x <- gsub(">", "$>$", x, fixed = TRUE)
18   x <- gsub("<", "$<$", x, fixed = TRUE)
19   x <- gsub("|", "$|$", x, fixed = TRUE)
20   x <- gsub("^", "\\verb|^|", x, fixed = TRUE)
21   x <- gsub("SANITIZE.BACKSLASH", "$\\backslash$",
22     x, fixed = TRUE)
23   x
24 }
25
26 ##' Retrieve column names required from QC
27 ##'
28 ##' @param datadir Directory name, ends in "/"
29 ##' @param fn file name
30 ##' @return vector of column names
31 ##' @author Paul Johnson
32 getQC <- function(datadir, fn = "QC_Required_Columns.xlsx"){
33   require(openxlsx)
34   qc <- read.xlsx(paste0(datadir, "/", fn), colNames = TRUE)
35   colnames(qc)
36 }
37
38 ## dts: delete trailing slash
39 ## Will delete repeat slashes first as well
40 dts <- function(name) gsub("/$", "", dms(name))
41
42 ##' Replace multiple slashes with one slash
43 dms <- function(name) gsub("(//)\1+", "/", name)
44
45
46 keywrds <- c("District", "County", "Middle.School", "High.School", "School",
47   "Elementary", "Northeast", "Southwest", "Northwest", "Southeast",
48   "Academy", "Cooperative", "Regional", "Foundations", "Preparatory",
49   "Community", "For.The.Blind.and.Visually.Impaired", "Independent",
50   "Special.Education",
51   "DISTRICT", "COUNTY", "SCHOOL", "ELEMENTARY", "MIDDLE.SCHOOL",
52   "HIGH.SCHOOL", "NORTHEAST", "SOUTHWEST", "NORTHWEST", "SOUTHEAST",
53   "ACADEMY", "COOPERATIVE", "REGIONAL", "FOUNDATIONS", "PREPARATORY",
54   "COMMUNITY", "FOR.THE.BLIND.AND.VISUALLY.IMPAIRED", "INDEPENDENT",
55   "SPECIAL.EDUCATION")
56
57
```

```

58 keyabr <- c("Dist", "Co", "MS", "HS", "Schl",
59           "Elem", "NE", "SW", "NW", "SE",
60           "Acdmy", "Co-op", "Reg", "Found", "Prep",
61           "Comm", "For_BVI", "Ind", "SPED",
62           "DIST", "CO", "SCHL", "ELEM", "MS",
63           "HS", "NE", "SW", "NW", "SE",
64           "ACDMY", "CO-OP", "REG", "FOUND", "PREP",
65           "COMM", "FOR_BVI", "IND", "SPED")
66
67 names(keywrds) <- keyabr
68 states.long <- c(state.name[48], state.name[-48], toupper(state.name)[48], toupper
        (state.name)[-48], "Choctaw", "Miccosukee")
69 states.short <- c(state.abb[48], state.abb[-48], state.abb[48], state.abb[-48], "
        Choctaw", "Miccosukee")
70 names(states.short) <- states.long
71
72 abbr <- function(x, long, short = NULL){
73   if (is.null(short)) short = names(long)
74   for (i in seq_along(long)){
75     x <- gsub(long[i], short[i], x)
76   }
77   x
78 }
79
80 shorten <- function(x, k = 20){
81   y <- substr(x, 1, k)
82   make.unique(y)
83 }
84
85 shrt <- function(x, k = 20) shorten(abbr(abbr(x, states.long, states.short),
        keywrds), k = 20)
86
87 abbrState1st <- function(x, state, long = states.long, short = states.short){
88   if (is.null(names(short))) stop("short name is unnamed")
89   xx <- ifelse(!is.null(short[state]), gsub(paste0("^", state), short[state], x
        ), x)
90   xx
91 }
92
93
94 ##' remove non-alpha numeric characters
95 ##'
96 ##' regular expression matching
97 ##' @param x
98 ##' @return
99 ##' @author Paul Johnson
100 clnstrng <- function(x){
101   y <- gsub("[^a-zA-Z0-9\\\\' _()-]", "", x)
102   y <- gsub("\\ *$", "", y) # Trailing spaces
103   y
104 }
105
106 ## clean file name string component
107 cfn <- function(x){
108   xn <- gsub("\\* $", "", x)
109   xn <- gsub("\\\\", "", x, fixed = TRUE)
110   xn <- gsub("\\ ", "_", xn)
111   xn <- gsub("_-", "_", xn)
112   xn <- gsub("_-", "_", xn)

```

```

113     xn <- gsub("-", "-", xn)
114     xn <- gsub("&", "and", xn)
115     xn <- gsub("\\\\'", "", xn)
116     xn <- gsub("\\\\\"", "", xn)
117     xn <- gsub("#", "-", xn, fixed = TRUE)
118     xn <- gsub("c/o", "co", xn)
119     xn <- gsub("/", "-", xn)
120     xn <- gsub("\\\\*", "", xn)
121     xn <- gsub("\\\\$", "", xn)
122     xn <- gsub("-$", "", xn)
123     xn <- gsub("_$", "", xn)
124     xn
125 }
126
127 ## recursively remove null objects
128 ## See: http://stackoverflow.com/questions/26539441/r-remove-null-elements-from-
129     list-of-lists
129 is.NullObs <- function(x) is.null(x) | all(sapply(x, is.null))
130
131 ## Recursively step down into list, removing all such objects
132 rmNullObs <- function(x) {
133     x <- Filter(Negate(is.NullObs), x)
134     lapply(x, function(x) if (is.list(x)) rmNullObs(x) else x)
135 }
136
137
138
139
140 ##' Remove rows that are more than 90% NA
141 ##'
142 ##' Some imported Excel sheets have "noise" in the last
143 ##' rows, which appear as <NA> or NA in R. Get rid of them.
144 ##' @title
145 ##' @param dframe
146 ##' @return
147 ##' @author Paul Johnson
148 deleteBogusRows <- function (dframe){
149     rowna <- apply(dframe, 1, function(x){sum(is.na(x))})
150     badrows <- rowna > .9 * dim(dframe)[2]
151     if (any(badrows)){
152         cat(paste("deleteBogusRows Diagnostic\n"))
153         cat(paste("These rows from the data frame: ", deparse(substitute(dframe)),
154             "\n are being purged:"))
154         print(rownames(dframe)[badrows])
155         cat(paste("The bad content was\n"))
156         print(dframe[badrows, ])
157         newdframe <- dframe[!badrows,]
158         return(newdframe)
159     }
160     print(paste("No bogus rows were found in: ", deparse(substitute(dframe))))
161     invisible(dframe)
162 }
163
164
165
166
167 ##' Searches a data frame to find column name components that match at least one
168 ##' element in a specified set of characters

```

```

169 ##'
170 ##' @param colnames A vector of words or symbols to use when searching
171 ##'     for matches
172 ##' @param data A data frame with named columns to be searched
173 ##' @param lower Logical. Determines whether or not the tolower()
174 ##'     function is used to ensure that all column names and search
175 ##'     elements match in terms of case
176 ##' @param names Logical. Determines whether or not the names of the
177 ##'     columns are returned by the function. The default is to
178 ##'     simply return the column numbers.
179 ##' @return Returns either the names of the columns, or the column
180 ##'     numbers that match at least one element in colnames.
181 ##' @author Ben Kite
182 cs <- function(strngs, data, lower = TRUE, names = TRUE){
183   names.orig <- names(data)
184   if (lower == TRUE){
185     varlist <- tolower(names.orig)
186     strngs <- tolower(strngs)
187   }else{
188     varlist <- names.orig
189     strngs <- strngs
190   }
191
192   varlist.orig <- varlist
193   ## whittle down varlist candidates iteratively (lexicographic search)
194   for (i in strngs){
195     varlist <- varlist[grep(i, varlist)]
196   }
197
198   if (length(varlist) > 1) stop(paste("There is more than one candidate in cs",
199     varlist))
200
201   winner <- match(varlist, varlist.orig)
202
203   if (names) return(names.orig[winner]) else winner
204 }
205
206 ##' Insert 0's in the front of existing digits or characters so that
207 ##' all elements of a vector have the same number of characters.
208 ##'
209 ##' The main purpose was to correct ID numbers in studies that are
210 ##' entered as integers with leading 0's as in 000001 or 034554. R's
211 ##' default coercion of integers will throw away the preceding 0's,
212 ##' and reduce that to 1 or 34554. The user might want to re-insert
213 ##' the 0's, thereby creating a character vector with values "000001"
214 ##' and "045665".
215 ##'
216 ##' If x is an integer or a character vector, the result is the
217 ##' more-or-less expected outcome (see examples). If x is numeric,
218 ##' but not an integer, then x will be rounded to the lowest integer.
219 ##'
220 ##' The previous versions of this function failed when there were
221 ##' missing values (NA) in the vector x. This version returns NA for
222 ##' those values.
223 ##'
224 ##' One of the surprises in this development was that sprintf() in R
225 ##' does not have a known consequence when applied to a character
226 ##' variable. It is platform-dependent (unpredictable). On Ubuntu Linux
227 ##' 16.04, for example \code{sprintf("%05s", 2)} gives back

```

```

227 ##' \code{" 2"}, rather than (what I expected) \code{"00002"}. The
228 ##' problem is mentioned in the documentation for \code{sprintf}. The
229 ##' examples show this does work now, but please test your results.
230 ##' @param x vector to be converted to a character variable by
231 ##'     inserting 0's at the front. Should be integer or character,
232 ##'     floating point numbers will be rounded down. All other
233 ##'     variable types will return errors.
234 ##' @param n Optional parameter. The desired final length of
235 ##'     character vector. This parameter is a guideline to determine
236 ##'     how many 0's must be inserted. This will be ignored if
237 ##'     \code{n} is smaller than the number of characters in the
238 ##'     longest element of \code{x}.
239 ##' @return A character vector
240 ##' @author Paul Johnson <pauljohn@ku.edu>
241 ##' @examples
242 ##' x1 <- c(0001, 0022, 3432, NA)
243 ##' padW0(x1)
244 ##' padW0(x1, n = 10)
245 ##' x2 <- c("1", "22", "3432", NA)
246 ##' padW0(x2)
247 ##' ## There's been a problem with strings, but this works now.
248 ##' ## It even preserves non-leading spaces. Hope that's what you want.
249 ##' x3 <- c("1", "22 4", "34323 42", NA)
250 ##' padW0(x3)
251 ##' x4 <- c(1.1, 334.52, NA)
252 ##' padW0(x4)
253 padW0 <- function (x, n = 0) {
254   if (!is.numeric(x) && !is.character(x))
255     stop("padW0 only accepts integer or character values for x")
256
257   xismiss <- is.na(x)
258   ## sprintf trouble with characters, not platform independent
259   if(is.numeric(x)) {
260     if (is.double(x)) x <- as.integer(x)
261     maxlength <- max(nchar(x), n, na.rm = TRUE)
262     vtype <- "d"
263     res <- sprintf(paste0("%0", maxlength, vtype), x)
264   } else {
265     maxlength <- max(nchar(x), n, na.rm = TRUE)
266     xlength <- vapply(x, nchar, integer(1))
267     padcount <- maxlength - xlength
268     pads <- vapply(padcount, function(i){paste0(rep("0", max(0, i, na.rm=TRUE)
269       ), collapse = "")}, character(1))
270     res <- paste0(pads, x)
271   }
272   res[xismiss] <- NA
273   res
274 }
275
276
277
278
279
280 ##' Cuts a string at a specified linewidth, trying to align cut with a
281 ##' separator
282 ##'
283 ##' Some strings are simply too long. Let's shorten them and attempt
284 ##' to cut at a specified separator. Use vectorized function,

```

```

285 ##' shortenv, to apply to non-scalar
286 ##' @param textstring target string to shorten
287 ##' @param linewidth approximate length to shorten taking into account
288 ##' tolerance
289 ##' @param tol number of characters forward/back to check; if single
290 ##' value then only backwards checking
291 ##' @param capwidth integer specifying the width of capital letters
292 ##' @param separator accepts 1- or 2-element vector to separate
293 ##' string; default = c(" ", "_")
294 ##' @return shortened character string
295 ##' @author Brent Kaplan, Ben Kite, Paul Johnson
296 ##' @examples
297 ##' test <- "123_567_9"
298 ##' (truncsmart(test, 5, tol = c(1,2)))
299 ##'
300 truncsmart <- function(textstring, linewidth, tol = c(5, 5),
301                          capwidth = 1.2, separator = c(" ", "_")) {
302   if(length(tol) > 2) stop("Please specify 1 or 2 values for tol.")
303   lets <- unlist(strsplit(textstring, split = ""))
304   caps <- sapply(lets, function(x) x %in% LETTERS)
305   widthval <- ifelse(caps == TRUE, capwidth, 1)
306   linelength <- cumsum(widthval)
307   if(linelength[length(linelength)] <= linewidth) return(textstring)
308   withinlength <- linelength <= linewidth
309   ifelse(length(tol) == 2, withintol <- linelength <= (linewidth + tol[2]),
310         withintol <- linelength <= (linewidth + tol))
311   letswl <- lets[withinlength]
312   letso <- letswl
313   letswt <- lets[withintol]
314   if(length(tol) == 2) {
315     for (j in 0:tol[1]){
316       if(identical(letswl[length(letswl)], separator[1]) ||
317          identical(letswl[length(letswl)], separator[2])) {
318         letso1 <- letswl[-length(letswl)]
319         break()
320       }
321       letswl <- letswl[-length(letswl)]
322     }
323     for (j in 0:tol[2]) {
324       if(identical(letswt[length(letswt)], separator[1]) ||
325          identical(letswt[length(letswt)], separator[2])) {
326         letso2 <- letswt[-length(letswt)]
327         break()
328       }
329       letswt <- letswt[-length(letswt)]
330     }
331     ifelse(!exists("letso1") && !exists("letso2"), letso,
332           ifelse(exists("letso1") && !exists("letso2"), letso <- letso1,
333                 ifelse(!exists("letso1") && exists("letso2"), letso <- letso2,
334                       ifelse(abs((length(letso1) - length(lets[withinlength]))) <= (length(
335                             letso2) - length(lets[withintol])), letso <- letso1, letso <- letso2)))
335   } else {
336     for (j in 0:tol){
337       if(identical(letswl[length(letswl)], separator[1]) ||
338          identical(letswl[length(letswl)], separator[2])) {
339         letso <- letswl[-length(letswl)]
340         break()
341       }

```

```

342     letswl <- letswl[-length(letswl)]
343   }
344 }
345 out <- paste0(letso, collapse = "")
346 out
347 }
348 truncsmart <- Vectorize(truncsmart, USE.NAMES=FALSE)
349
350
351
352 ##' Initial letter will be capitalized
353 ##'
354 ##' copied from Frank Harrell's Hmisc package
355 ##' @param string argument to be capitalized
356 ##' @return A capitalized string
357 ##' @author Frank Harrell
358 capitalize <- function (string)
359 {
360   capped <- grep("[^A-Z]*$", string, perl = TRUE)
361   substr(string[capped], 1, 1) <- toupper(substr(string[capped],
362     1, 1))
363   return(string)
364 }
365
366
367
368 ##' Cuts a string at a specified linewidth, trying to align cut with a separator
369 ##'
370 ##' Some strings are simply too long. Let's shorten them and attempt
371 ##' to cut at a specified separator
372 ##' @param textstring
373 ##' @param linewidth
374 ##' @param separator defaults to space but can be underscore "_"
375 ##' @return shortened character string
376 ##' @author Ben Kite, Paul Johnson, Brent Kaplan
377 ##' @examples
378 ##' test <- "An example string that is way too long to be useful and therefore
379 ##' needs to be cut down to our specified linewidth."
380 ##'
381 splitter <- function(textstring, linewidth, separator = " "){
382   lets <- unlist(strsplit(textstring, split = ""))
383   caps <- sapply(lets, function(x) x %in% LETTERS)
384   widthval <- ifelse(caps == TRUE, 1.2, 1)
385   linelength <- cumsum(widthval)
386   if(linelength[length(linelength)] < linewidth) return(textstring)
387   withinlength <- linelength < linewidth
388   lets2 <- lets[withinlength]
389   lets3 <- lets2
390   for (j in 1:5){
391     if(lets3[length(lets3)] == separator) {
392       lets3 <- lets3[-length(lets3)]
393       lets2 <- lets3
394       break()
395     }
396     lets3 <- lets3[-length(lets3)]
397   }
398   out <- paste0(lets2, collapse = "")
399   out

```

```

400 }
401
402
403 ##' Replace column names with new names from a named vector
404 ##'
405 ##' Imagine some beautiful writing here
406 ##' @param dat a data frame
407 ##' @param newname A named vector of the form c(oldname1 = "newname1", oldname2 =
      "newname")
408 ##' @param lowercase should new names be converted to lower case?
409 ##' @return a data frame
410 ##' @export
411 ##' @author Paul Johnson <pauljohn@ku.edu>
412 colnamesReplace <- function(dat, newname, lowercase = FALSE){
413   if (is.null(names(newname))){
414     stop("The newname vector should be a named vector, using the oldnames as
          the names")
415   }
416   oldname <- colnames(dat)
417   if(any(newname[oldname] != oldname)){
418     cat(paste("colnamesReplace Diagnostic: These names are to be replaced:\n")
        )
419     res <- ifelse (oldname %in% names(newname), newname[oldname], oldname)
420     print(c("oldname", "newname"))
421     print(cbind(oldname, res))
422     colnames(dat) <- if (lowercase) tolower(res) else res
423     return(dat)
424   } else{
425     cat(paste("replaceNames Diagnostic: No column names were altered"))
426     return(dat)
427   }
428 }

```