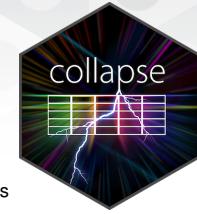


# Advanced and Fast Data Transformation with `collapse` :: CHEAT SHEET



## Introduction

`collapse` is a C/C++ based package supporting advanced (grouped, weighted, time series, panel data and recursive) statistical operations in R, with very efficient low-level vectorizations across both groups and columns.

It also offers a flexible, class-agnostic, approach to data transformation in R: handling matrix and data frame based objects in a uniform, attribute preserving, way, and ensuring seamless compatibility with `dplyr` / (grouped) `tibble`, `data.table`, `xts`, `sf` and `plm` classes for panel data ('`pseries`', '`pdata.frame`').

`collapse` provides full control to the user for statistical programming - with several ways to reach the same outcome and rich optimization possibilities. Its default is `na.rm = TRUE`, and implemented at very low cost at the algorithm level.

Calling `help("collapse-documentation")` brings up a detailed documentation, which is also available [online](#). See also the `fastverse` package/project for a recommended set of complimentary packages and easy package management.

## Row/Column Arithmetic (by Reference)

Column-wise sweeping out of vectors/matrices/DFs/lists

```
%cr%, %c+, %c-, %c*, %c/% e.g. Z = X %c% rowSums(X)
```

Row-wise sweeping vectors from vectors/matrices/DFs/lists

```
%rr%, %r+, %r-, %r*, %r/% e.g. Z = X %r% colSums(X)
```

Standard (column-wise) math by reference (returns invisibly)

```
%+=%, %-=%, %*=%, %/=% e.g. X %-=% rowSums(X)
```

Same thing, also supports row-wise operations by reference

```
setop(X, "/", rowSums(X))
setop(X, "/", colSums(X), rowwise = TRUE)
```

## Transform Data by (Grouped) Replacing or Sweeping out Statistics (by Reference)

A generalisation of rowwise operations, that also supports sweeping by groups e.g. aggregate statistics

```
TRA(x, STATS, FUN = "-", g = NULL, set = FALSE)
setTRA(x, STATS, FUN = "-", g = NULL)
```

x vector, matrix, or (grouped) data frame / list

STATS statistics matching (columns of) x (i.e. aggregated vector, matrix or data frame / list)

FUN integer/string indicating transformation to perform:

Int.	String	Description
0	"replace_NA"	replace missing values in x
1	"replace_fill"	replace data and missing values in x
2	"replace"	replace data but preserve missing values in x
3	"+"	subtract: x - STATS(g)
4	"-"	x - STATS(g) + fmean(STATS, w = GRP)
5	"*"	divide: x / STATS(g)
6	"%"	compute percentages: x * 100/STATS(g)
7	"*%"	add: x + STATS(g)
8	"%"	multiply: x * STATS(g)
9	"%%"	modulus: x %% STATS(g)
10	"-%%"	subtract modulus: x - x %% STATS(g)

g [optional] (list of) vectors / factors or GRP() object

set TRUE transforms x by reference. setTRA is equivalent to invisible(TRA(..., set = TRUE))

## Fast Statistical Functions

Fast functions to perform column-wise grouped and weighted computations on matrix-like objects

```
fmean, fmedian, fmode, fsum, fprod, fsd, fvar
fmin, fmax, fnth, ffist, flast, fnobs, fndistinct
```

## Syntax

```
FUN(x, g = NULL, [w = NULL], TRA = NULL,
[na.rm = TRUE], use.g.names = TRUE,
[drop = TRUE], [nthreads = 1L])

x vector, matrix, or (grouped) data frame / list
g [optional] (list of) vectors / factors or GRP() object
w [optional] vector of (frequency) weights

TRA [optional] operation to transform data with computed
statistics (see FUN argument to TRA() and Examples)

drop drop matrix / data frame dimensions. default TRUE
```

## Examples

```
fmean(AirPassengers) # Vector
## [1] 280.2986
fmean(AirPassengers, w = cycle(AirPassengers)) # Weighted mean
## [1] 284.3397
fmean(EuStockMarkets) # Matrix
##   DAX    SMI    CAC    FTSE
## 2530.657 3376.224 2227.828 3565.643
fmean(EuStockMarkets, drop = FALSE) # Don't drop dimensions
##   DAX    SMI    CAC    FTSE
## [1,] 2530.657 3376.224 2227.828 3565.643
fmean(airquality) # Data Frame (can also use drop = FALSE)
##   Ozone Solar.R  Wind   Temp Month Day
## 42.129310 185.931507  9.957516 77.882353 6.993464 15.803922
fmean(iris[1:4], g = iris$Species) # Grouped
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## setosa      5.006      3.428      1.462      0.246
## versicolor  5.936      2.770      4.260      1.326
## virginica   5.858      2.974      5.525      2.026
X = iris[1:4]; g = iris$Species; w <- abs(rnorm(nrow(X)))
fmean(X, g, w) # Grouped and weighted (random weights)
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## setosa      5.011663  3.467638  1.504067  0.2525002
## versicolor  5.930365  2.773558  4.238593  1.3136082
## virginica   5.88903  2.978017  5.552375  2.0221178
## Transformations: here center data on the weighted group median
TRA(X, fmedian(X, g, w), "-", g) |> head(3)
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1          0.1        0.0       -0.1        0
## 2         -0.1       -0.5       -0.1        0
## 3         -0.3       -0.3       -0.2        0
fmedian(X, g, w, TRA = "-") |> head(3) # Same thing: more compact
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1          0.1        0.0       -0.1        0
## 2         -0.1       -0.5       -0.1        0
## 3         -0.3       -0.3       -0.2        0
fmedian(X, g, w, "-", set = TRUE) # Modify in-place (same as setTRA())
head(iris, 3) # Changed iris too, as X = iris[1:4] did a shallow copy
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          0.1        0.0       -0.1        0     setosa
## 2         -0.1       -0.5       -0.1        0     setosa
## 3         -0.3       -0.3       -0.2        0     setosa
```

## Basic Computing with R Functions

Apply R functions to rows or columns (by groups)

```
dapply(x, FUN, ..., MARGIN = 2) - column/row apply
BY(x, g, FUN, ...) - split-apply-combine computing
```

## Grouping and Ordering

Optimized functions for grouping, ordering, unique values, splitting & recombinining, and dealing with factors

GRP() - create a grouping object (class 'GRP'): pass to g arg.

```
g <- GRP(iris, ~ Species) # or GRP(iris$Species) or GRP(iris[["Species"]])
fndistinct(iris[1:4], g) # Computation without grouping overhead
```

```
##           Sepal.Length Sepal.Width Petal.Length Petal.Width
## setosa            15          16          9          6
## versicolor        21          14          19          9
## virginica         21          13          20          12
```

fgroup\_by() - attach 'GRP' object to data: a class-agnostic grouped frame supporting fast computations

```
mtcars |> fgroup_by(cyl, vs, am) |> ss(1:2)
```

```
##           mpg cyl disp hp drat wt qsec vs am gear carb
## Mazda RX4   21   6 160 110 3.9 2.620 16.46 0 1 4 4
## Mazda RX4 Wag 21   6 160 110 3.9 2.875 17.02 0 1 4 4
```

```
## Grouped by: cyl, vs, am [7 1 5 (3.8) 1-12]
```

```
# Group Stats: [N, groups] | mean (sd) min-max of group sizes]
# Fast Functions also have a grouped_df method: here ut-weighted medians
```

```
mtcars |> fgroup_by(cyl, vs, am) |> fmedian(trt) |> head(3)
```

```
##   cyl vs am sum.wt mpg disp hp drat qsec gear carb
## 1   4   0   1  2.140 26.0 120.3 91 4.43 16.70 5   2
## 2   4   1   0  8.805 22.8 140.8 95 3.70 20.01 4   2
## 3   4   1   1 14.198 30.4 79.0 6 4.06 18.61 4   1
```

GRPN(), fgrounp\_vars(), fungroup() - get group count, grouping columns/variables, and ungroup data

qF(), qG() - quick as.factor, and vector grouping object of class 'qG': a factor-light without levels attribute

group() - (multivariate) group id ('qG') in appearance order

groupid() - run-length-type group id ('qG')

seqid() - group-id from integer-sequences ('qG')

radixorder[v]() - (multivariate) radix-based ordering

finteraction() - fast factor interactions (or return 'qG')

fdroplevels() - fast removal of unused factor levels

f[n]unique() - fast unique values / rows (by columns)

gsplit() - fast splitting vector based on 'GRP' objects

greorder() - efficiently reorder y = unlist(gsplit(x, g)) such that identical(greorder(y, g), x)

collapse optimizes grouping using both factors / 'qG' objects and 'GRP' objects. 'GRP' objects contain most information and are thus most efficient for complex computations.

```
X <- iris[1:4]; v <- as.character(iris$Species)
```

```
f <- qF(v, na.exclude = FALSE) # Adds 'na.included' class: no NA checks
gv <- group(v) # 'qG' object: first appearance order, with 'na.included'
```

```
microbenchmark(fmode(X, v), fmode(X, f), fmode(X, gv), fmode(X, g))
```

```
# Unit: microseconds
```

```
##   expr      min       q1       median      max neval
## 1 fmode(X, v) 11.890 12.9150 13.17697 13.3455 100
## 2 fmode(X, f)  9.225  9.8195 11.33035 10.0860 10.4550 100
## 3 fmode(X, gv) 8.569  9.3480 10.73667 9.6555 10.1065 100
## 4 fmode(X, g)  6.683  7.2980  7.71620 7.5440  7.7490 100
```

```
# Population weighted mean (PCGD, LIFEEX) & mode (country), and sum(POP)
```

```
collap(wldeid, country + PCGD + LIFEEX ~ income, w = "POP")
```

```
##   country      income      PCGD      LIFEEX      POP
## 1 United States High income 31284.7366 75.69257 58840837058
```

```
## 2 Ethiopia Low income 557.1427 53.50608 20949161394
```

```
## 3 India Lower middle income 1238.8280 60.58651 113837684528
```

```
## 4 China Upper middle income 4145.6844 68.26984 119606023798
```

## Advanced Transformations

Common transformations (in econometrics)

Scaling, Centering and Averaging

```
fscale(x, g = NULL, w = NULL, na.rm = TRUE,
       mean = 0, sd = 1, ...)
fwithin(x, g = NULL, w = NULL, na.rm = TRUE,
        mean = 0, theta = 1, ...)
fbetween(x, g = NULL, w = NULL, na.rm = TRUE,
          fill = FALSE, ...)
```

Higher-Dimensional Centering/Avg. and Linear Prediction

```
fhdwithin(x, fl, w = NULL, na.rm = TRUE,
           fill = FALSE, lm.method = "qr", ...)
fhdbetween() - same arguments as fhdwithin()
```

Statistical Operators (function shorthands with extra features)

```
STD(), W(), B(), HDW(), HDB()
```

## Examples

```
# Grouped scaling
iris > fgroup_by(Species) |> fscale() |> head(2)
## Species Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1 setosa   0.2666745  0.1899414 -0.3570112 -0.4364923
## 2 setosa   0.2666745  0.1899414 -0.3570112 -0.4364923
STD(iris, ~Species, stub = FALSE) # Same thing + faster
# Grouped and weighted scaling. Operators support formulas and keep ids
STD(mtcars, mpg ~ carb ~ cyl, w = "wt") |> head(2)
##      cyl   wt STD.mpg STD.carb
## Mazda RX4 6 2.620 0.9691687 0.386125
## Mazda RX4 Wag 6 2.875 0.9691687 0.386125
# Much shorter than fssubset(mpg > fmean(mpg, cyl, TRA = "replace"))
mtcars > fssubset(mpg ~ B(cyl, mpg)) |> head(2)
##      mpg cyl disp hp drat wt qsec vs am gear carb
## Mazda RX4 21  6 160 110 3.9 2.620 16.46 0 1 4 4
## Mazda RX4 Wag 21  6 160 110 3.9 2.875 17.02 0 1 4 4
# Regression with cyl fixed effects - a la Mundlak (1978)
lm(mpg ~ carb + B(carb, cyl), data = mtcars) |> coef()
## (Intercept) carb B(carb, cyl)
## 34.829652 -0.465511 -4.775032
# Fast grouped (vs) bivariate regression slopes: mpg ~ carb
mtcars > fgroup_by(vs) |> fmutate(dm_carb = W(carb)) |>
  summarise(beta = fsum(mpg, dm_carb) %~% fsum(dm_carb^2))
## vs      beta
## 1 0 -0.5557241
## 2 1 -2.0706468
# Residuals from regressing on 'Petal' vars and 'Species' FE
fhdwithin(iris[1:2], iris[3:5]) |> head(2)
## Sepal.Length Sepal.Width
## 1 0.14989286 0.1102684
## 2 0.20501074 -0.3897316
# Detrending with country-level cubic polynomials
HDW(wlddev, PCGDP + LIFEEX + POP ~ iso3c * poly(year, 3)) |> head(2)
## HDW.PCGDP HDW.LIFEEX HDW.POP
## 43 -258.4069 0.2360285 -317459.1
## 44 -119.5600 0.1136432 -33900.2
# Note: HD centering/prediction and polynomials requires package 'fjtest'
```

## Linear Models

Fast (barebones) linear model fitting with 6 different solvers

```
flm(y, X, w = NULL, add.icpt = FALSE, method = "lm")
```

Fast  $R^2$ -based F-test of exclusion restrictions for lm's (with FE)

```
fFtest(y, exc, X = NULL, w = NULL, full.df = TRUE)
```

Both functions also have formula interfaces:

```
flm(cbind(mpg, disp) ~ hp + carb, weights = wt, mtcars)
##      mpg      disp
## (Intercept) 28.48041839 42.155002
## hp         -0.06834996  2.101036
## carb        0.33207257 -38.183910
# Test the exclusion of cyl-dummies and hp.
fFtest(mpg ~ qB(cyl) + hp | carb + qF(am), weights = wt, mtcars)
##      R-Sq. DF1 DF2 F-Stat. P-Value
## Full Model    0.812  5  26 22.479  0.000
## Restricted Model 0.674  2  29 30.041  0.000
## Exclusion Rest. 0.138  3  26  6.351  0.002
```

## Time Series and Panel Series

Fast and flexible indexed series and data frames: a modern upgrade of plm's 'pseries' and 'pdata.frame'

Turn DF into an 'indexed\_frame' using id and/or time vars

```
data_ix = findex_by(data, id1, ..., time)

data_ix$indexed_series - columns are 'indexed_series'
index_df = findex(data_ix) - retrieve 'index_df': DF of ids
index_df = with(data_ix, findex(indexed_series)) - can fetch 'index_df' from 'indexed_series' in any caller environment
data = unindex(data_ix) - unindex (also 'indexed_series')
reindex(data, index = index_df) - reindex / new pointers
'indexed_series' can be 1-or-2D atomic objects. Vectors / time series / matrices can also be indexed directly using:
reindex(vec/mat, index = vec/index_df)

is_irregular() - irregularity in any index[ed] obj. or time vec
```

### Example: Indexing Panel Data

```
wldi <- wlddev |> findex_by(iso3c, year) # Balanced: 216 countries
fssubset(wldi, 1:12, iso3c, year, PCGDP:POP)
## iso3c year PCGDP LIFEEX GINI ODA POP
## 1 AFG 1960 NA 32.446 116769997 8996973
## 2 AFG 1961 NA 32.962 232080002 9169410
##
## Indexed by: iso3c [1] | year [2 (61)]
## Index stats: [N. ids] / [N. periods (tot.N. periods: (max-min)/GCD)]
LIFEEXI = wldi$LIFEEX # Indexed series
str(LIFEEXI, strict.width = "cut")

## 'indexed_series' num [1:13176] 32.4 33 33.5 34 34.5 ...
## - attr(*, "index_df")=Classes 'index_df', 'pindex' and 'data.frame'...
## .. $ iso3c: Factor w/ 216 levels "ABW","AFG","AGO",... 2 2 2 2 2 2 ...
## .. $ year: Ord.factor w/ 61 levels "1960"<-"1961"<...: 1 2 3 4 5 6 ...
LIFEEXI[1:7] # Subsetting indexed series
```

Note: 'indexed\_series' and frames are supported via existing 'pseries'/pdata.frame' methods for time series/panel functions.

Fast functions to perform time-based computations on (irregular) time series and (unbalanced) panel data

Lags/Leads, Differences, Growth Rates and Cumulative Sums

```
flag(x, n = 1, g = NULL, t = NULL, fill = NA, ...)
fdiff(x, n = 1, diff = 1, g = NULL, t = NULL,
      fill = NA, log = FALSE, rho = 1, ...)
fgrowth(x, n = 1, diff = 1, g = NULL, t = NULL, fill = NA, logdiff = FALSE, scale = 100, power = 1, ...)
fcumsum(x, g = NULL, o = NULL, na.rm = TRUE, fill = FALSE, check.o = TRUE, ...)
```

Statistical Operators: L(), F(), D(), Dlog(), G()

### Example: Computing Growth Rates

```
# Ad-hoc use: note that G() supports formulas which fgrowth() doesn't
fgrowth(AirPassengers) |> head()
## [1] NA 5.357143 11.864407 -2.272727 -6.201550 11.570248
G(wlddev, c(1, 10), by = PCGDP ~ iso3c, t = ~year) |> ss(11:12)
## iso3c year G1.PCGDP L10G1.PCGDP
## 1 AFG 1970 NA NA
## 2 AFG 1971 NA NA
wlddev |> fgroup_by(iso3c) |> fselect(iso3c, year, PCGDP, LIFEEX) |>
  fmutate(PCGDP_growth = fgrowth(PCGDP, t = year)) |> head(2)
## iso3c year PCGDP LIFEEX PCGDP_growth
## 1 AFG 1960 NA 32.446 NA
## 2 AFG 1961 NA 32.962 NA
settransform(wlddev, PCGDP_growth = G(PCGDP, g = iso3c, t = year))
# Note: can omit t -> requires consecutive observations and groups
# Usage with indexed series / frames:
```

## Recode and Replace Values

```
recode_num(), recode_char() - recode numeric / character values (+ regex recoding) in matrix-like objects
replace_[NA|Inf|outliers]() - replace special values
pad() - add (missing) observations / rows i.e. expand objects
```

## (Memory) Efficient Programming

Functions for (memory) efficient R programming

```
any/all[v|NA], which[v|NA], %[=!=]=%, copyv, setv, alloc
missing_cases, na_[insert|rm|omit], vlengths, vtypes,
vgcd, frange, fnlevels, fn[row|col], fdim, seq_[row|col]
fssubset(wlddev, year %~% 2010) # 2x faster fssubset(wlddev, year == 2010)
attach(mtcars) # Efficient sub-assignment by reference, various options...
setv(am, 0, vs); setv(am, 1:10, vs); setv(am, 1:10, vs[10:20])
```

## Small (Helper) Functions

Functions for (meta-)programming and attributes

```
.c, massign, %=%, vlabels[<-], setLabels, vclasses,
namlab, [add]rm_stub, %!in%, ckmatch, all_identical,
all_obj_equal, all_funs, set[Dim|Row|Col]names,
unattrib, setattr, copyAttrib, copyMostAttrib
.c(var1, var2, var3) # Non-standard concatenation
## [1] "var1" "var2" "var3"
c(values, vectors) %~% eigen(cov(mtcars)) # Multiple Assignment
# Variable labels: vlabels[<-], isetjrelabel() etc. namlab() shows summary
namlab(wlddev[c(2, 9)], N = TRUE, Ndist = TRUE, class = TRUE)
## Variable Class N Ndist Label
## 1 iso3c factor 13176 216 Country Code
## 2 PCGDP numeric 9470 9470 GDP per capita (constant 2010 US$)
```

## API Extensions

Shorthands for frequently used functions

```
fselect -> slt, fssubset -> sbt, fmutate -> mtt,
[fc|set]transform[v] -> [set]tfm[v], fsummarise -> smr,
across -> acr, fgroup_by -> gby, finteracter -> itn,
findex_by -> iby, findex -> ix, frename -> rnm,
get_vars -> gv, num_vars -> nv, add_vars -> av
```

## Namespace masking

Can set option(collapse\_mask = c(...)) with a vector of functions starting with f, to export versions without f-, masking base R or dplyr. A few keywords exist to mask multiple functions, see help("collapse-options"). This allows clean & fast code, but poses additional namespace challenges:

```
# Masking all f- functions and specials n = GRPN and table = qtab
options(collapse_mask = "all")
library(collapse)
# The following is 100% collapse code, apart from the base pipe
```

```
wlddev |>
  subset(year >= 1990) |>
  group_by(year) |>
  summarise(n = n(), across(PCGDP:GINI, mean, w = POP))
```

```
with(mtcars, table(cyl, vs, am))
sum(mtcars)
diff(EuStockMarkets)
dropLevels(wlddev)
mean(nv(iris), g = iris$Species)
scale(nv(GGDC10S), g = GGDC10S$variable)
unique(GGDC10S, cols = c("Variable", "Country"))
range(wlddev$date)
```

```
wlddev |>
  index_by(iso3c, year) |>
  mutate(PCGDP_lag = lag(PCGDP),
         PCGDP_diff = PCGDP - PCGDP_lag,
         PCGDP_growth = growth(PCGDP)) |> unindex()
```

The best way to set this option is inside an .Rprofile file placed in the user or project directory. Use it carefully.