

Why the package **declared**

The problem

The R ecosystem already has some very good packages that deal with labelled objects. In particular, the inter-connected packages **haven** and **labelled** provide all the functionality most users would ever need.

As nice and useful as these packages are, it has become apparent they have some fundamental design features that run, in some situations, against users' expectations. This has a lot to do with the treatment of declared missing values, that are instrumental for the social sciences.

The following minimal example (adapted from the vignette in package **haven**) illustrates the situation:

```
> library(haven)
> x1 <- labelled_spss(c(1:5, 99), labels = c(Missing = 99), na_value = 99)
```

The printed objects from this package nicely display some properties:

```
> x1
<labelled_spss<double>[6]>
[1] 1 2 3 4 5 99
Missing values: 99

Labels:
  value  label
    99 Missing
```

There are 5 normal (non-missing) values (supposedly they represent the number of children), and one declared missing value coded 99. This value *acts* as a missing value, but it is different from a regular missing value in R, coded NA. The latter stands for any missing information (something like an empty cell) regardless of the reason.

Here, on the other hand, the cell is *not* empty, but the value 99 is not a valid value either. It cannot possibly represent 99 children in the household, but for instance it could mean the respondent did not want to respond. It is properly identified as missing, with:

```
> is.na(x1)
[1] FALSE FALSE FALSE FALSE FALSE  TRUE
```

But when calculating a mean, for instance, the normal expectation is that value 99 would not play any role in the calculations (since it should be *missing*). However:

```
> mean(x1)
[1] 19
```

This means the value 99 did play an active role despite being identified as “missing”. In an ideal world, the expected mean would be 3, or at best employ the argument `na.rm = TRUE` if the result is NA because of the declared missing value.

A solution to this problem is offered by package **labelled**, which has a function called `user_na_to_na()`:

```
> library(labelled)
> mean(user_na_to_na(x1), na.rm = TRUE)
[1] 3
```

The declared solution

While solving the problem, this above solution forces two additional operations:

- converting the (already) declared user missing values, and
- employing the `na.rm` argument.

This should not be necessary, especially if (and it is extremely likely that) users may forget the declared missing values are not actually missing values. This scenario is quite possible, as many users previously using other software like SPSS or Stata where nothing else should be done after declaring the missing values, may not realize more is needed.

To solve this situation, package **declared** creates a very similar object, where declared missing values are actually stored (hence interpreted as) regular NA missing values in R.

```
> library(declared)
> x2 <- declared(c(1:5, 99), labels = c(Missing = 99), na_value = 99)
> x2
<declared<integer>[6]>
[1]      1      2      3      4      5 NA(99)
Missing values: 99

Labels:
  value  label
    99 Missing
```

It is now obvious the value 99 is not a regular number anymore, but an actual missing value. More importantly, it circumvents the need to convert user missing values to regular NAs, since they are already stored as NA values. The average value is calculated simply as:

```
> mean(x2)
[1] 3
```

Notice that neither `user_na_to_na()`, nor employing `na.rm = TRUE` are necessary and, despite being stored as an NA value, the value 99 is not equivalent to an *empty cell*. The information still exists, but it is simply ignored in the calculations.

The `na.rm = TRUE` is only necessary if there are truly unexplained (hence undeclared) missing values in the data:

```
> mean(c(x2, NA))
[1] NA
> mean(c(x2, NA), na.rm = TRUE)
[1] 3
```

As it can be seen, combining on this vector creates a similar one of the same class:

```
> x2 <- c(x2, 97, 99)
> x2
<declared<integer>[8]>
[1]      1      2      3      4      5 NA(99)    97 NA(99)
Missing values: 99

Labels:
  value  label
    99 Missing
```

Similarities and added value

It should be made obvious that packages **haven** and **labelled** are excellent packages which are not inherently doing a bad thing: the very same result is obtained, just via a different route. Package **declared** should not even be necessary, if the design philosophy of these packages would be different.

The current functions offer an alternative to these packages, with only one but fundamental difference: instead of treating existing values as missing, package **declared** interprets missing values as existing.

The declared missing values are (just like in package **haven**) identified as NAs, but they can also be compared against the original values:

```
> is.na(x2)
[1] FALSE FALSE FALSE FALSE FALSE  TRUE FALSE  TRUE
> x2 == 99
[1] FALSE FALSE FALSE FALSE FALSE  TRUE FALSE  TRUE
```

Most functions are designed to be as similar as possible, for instance `value_labels()` to add / change value labels:

```
> value_labels(x2) <- c("Do not know" = 97, "Not responded" = 99)
> x2
<declared<integer>[8]>
[1]      1      2      3      4      5 NA(99)    97 NA(99)
Missing values: 99

Labels:
  value      label
    97 Do not know
    99 Not responded
```

The value 97 is now properly labelled, and it can further be declared as missing. Such declarations do not necessarily have to use the main function `declared()`, due to the separate functions `missing_values()` and `missing_range()`:

```
> x2 <- c(x2, -3, -1, -2, -5)
> missing_values(x2) <- c(97, 99)
> missing_range(x2) <- c(-1, -5)
> x2
<declared<integer>[12]>
[1]      1      2      3      4      5 NA(99) NA(97) NA(99) NA(-3) NA(-1)
[11] NA(-2) NA(-5)
Missing values: 97, 99
Missing range:  [-5, -1]

Labels:
  value      label
    97 Do not know
    99 Not responded
```

To ease the smooth inter-operation with packages **haven** and **labelled**, the following functions are of interest: `undeclare()`, `as_haven()` and `as_declared()`.

The function `undeclare()` replaces the NAs with their declared missing values. The result is still an object of class `declared`, but all missing values (and missing range) are stripped off the vector and values are presented as they have been collected. All other attributes of interest (variable and value labels) are retained and printed accordingly.

The function `as_haven()` coerces the resulting object to the class `haven_labelled_spss`, and the function `as_declared()` reverses the process:

```

> xh <- as_haven(x2)
> xh
<labelled_spss<integer>[12]>
  [1]  1  2  3  4  5 99 97 99 -3 -1 -2 -5
Missing values: 97, 99
Missing range:  [-5, -1]

Labels:
  value      label
    97    Do not know
    99 Not responded
>
> as_declared(xh)
<declared<integer>[12]>
  [1]      1      2      3      4      5 NA(99) NA(97) NA(99) NA(-3) NA(-1)
[11] NA(-2) NA(-5)
Missing values: 97, 99
Missing range:  [-5, -1]

Labels:
  value      label
    97    Do not know
    99 Not responded

```

The declared objects play natively with the base functions `na.omit()` or `na.remove()`, either as standalone vectors or part of a data frame:

```

> dfm <- data.frame(id = sample(1:12, 12), x2)
> dfm
   id  x2
1  6    1
2  9    2
3  5    3
4  2    4
5  7    5
6 11 NA(99)
7  1 NA(97)
8 12 NA(99)
9  8 NA(-3)
10 4 NA(-1)
11 10 NA(-2)
12 3 NA(-5)

```

```

> na.omit(dfm)
  id x2
1  6  1
2  9  2
3  5  3
4  2  4
5  7  5

```

As it can be noticed, the missing values are properly formatted inside a data frame. Other users might prefer a tibble instead of a data frame, in which case objects of class **declared** are properly formatted in a similar way to those from package **haven**:

```

> library(tibble)
> dfm$xh <- as_haven(x2)
> as_tibble(dfm)
# A tibble: 12 x 3
   id          x2          xh
  <int>    <int+lbl>    <int+lbl>
1     6     1         1
2     9     2         2
3     5     3         3
4     2     4         4
5     7     5         5
6    11 99 (NA) [Not responded] 99 (NA) [Not responded]
7     1 97 (NA) [Do not know]  97 (NA) [Do not know]
8    12 99 (NA) [Not responded] 99 (NA) [Not responded]
9     8 -3 (NA)         -3 (NA)
10    4 -1 (NA)         -1 (NA)
11    10 -2 (NA)        -2 (NA)
12     3 -5 (NA)        -5 (NA)

```

There is an obvious amount of duplication between packages **declared** and **labelled**. As mentioned, the function `value_labels()` corresponds to the function `val_labels()` in package **labelled**, and the same overlap happens between functions `variable_labels()` and `var_labels()`.

The list could grow but it is quite unnecessary, as specific methods can be added to the functions from package **labelled**, to deal with the objects of class **declared**. And indeed, most of them work natively, for instance:

```

> x <- declared(c(1:5), labels = c(Good = 1, Bad = 5))
> to_factor(x)
[1] Good 2    3    4    Bad
Levels: Good 2 3 4 Bad

```

```
> to_character(x)
[1] "Good" "2"    "3"    "4"    "Bad"
```

There is no need to further duplicate these functions in package **declared**, as they work very well by simply loading the package **labelled**. Ideally, all of these packages should be merged to provide a unified experience, but otherwise package **declared** can be used independently (and gracefully coexist with the other two).