

R source code for loading LOS_model dataset from NHSRdatasets and its data dictionary

B213753

18/06/2022

Link to github repo

https://github.com/B213753/B213753_assessment.git

Introduction

The data capture tool will look at the (length of stay) 'LOS_model' dataset and attempt to model hospital length of stay (LOS) in the ten different trusts within the dataset. Being able to analyse the differences, if any, in length of stay in different trusts will allow for better planning and allocation of resources, or perhaps highlight areas for improvement with future more in-depth analyses.

The tool will look to collect the following variables: patient ID, organisation, patient age, length of stay, and death. Please see the data dictionary below for more details on these variables.

Script function and output

The purpose of this script is to load the LOS_model dataset into R, give a brief overview of the data using summary tables, recode the death variable into a factor and finally split the dataset into test and train sets.

The script also uses R to build a data dictionary that will be appended onto the data collected with the data capture tool.

The formal data dictionary is also provided in this document.

Outputs

The following files will be created/modified on successful conclusion of this script. This script assumes the working directory is set as the root project directory.

- */RawData/LOS_raw.csv* - The raw LOS data containing 300 observations.
- */Data/LOS_train.csv* - The train dataset containing 288 observations.
- */Data/LOS_test.csv* - The test dataset containing 11 observations.
- */Data/LOS_test_marker.csv* - The test marker dataset containing 1 observation.
- */RawData/CollectedData_DataDictionary.csv* - The data dictionary created in R saved as csv.
- */RawData/complete_CollectedData.rds* - R dataset file containing the data dictionary and other metadata appended onto the collected dataset.

Load packages and data

The following packages are required for this script:

- **NHSRdatasets** - the package containing the dataset selected (LOS_model) for this project.
- **tidyverse** - The tidyverse is a collection of R packages used in this script for most functions in this script in the manipulation of data and files.
- **here** - Used for easy file referencing in project-oriented workflows.
- **knitr** - Used in this script to visualize tables using kable().
- **caret** - Used in this script for a function called createDataPartition() which helped split the dataset into test and training sets.
- **dataMeta** - Used to create a data dictionary within R and append to a dataset.

```
library(NHSRdatasets)
library(tidyverse)
library(here)
library(knitr)
library(caret)
library(dataMeta)
```

Data Dictionary

See below for the R source code which builds and appends a dictionary in R to the collected dataset.

Data use

The data management plan is found in the Outputs folder of the repository. The data collection tool will collect the following variables for use in statistical analyses and modelling: patient ID, organisation, patient age, length of stay, and death.

The LOS_model is a synthetic dataset from the NHSRdatasets package in R. Data collected using the data collection tool is created during manual user inputs using the tool.

The dataset will undergo analysis, exploration and transformation using RStudio within the Author's Noteable environment. The dataset and its outputs will be used to analyse and model the length of stay (LOS) data for hospital admissions. The aim is to seek meaning relationships that aims to improve health and social care resource allocation and delivery.

Variable description

The collected dataset contains:

- **ID** : an integer value representing patient number - should be unique and allows referencing rows to the original raw data.
 - Quantitative data type.
 - Min value: none
 - Max value: 99999
 - Allowed values: only integer values allowed. The tool does not check for unique IDs (allows data entry for duplicate ID).
 - Example: 1, 43
- **Organisation** : a factor consisting of a string that represents a hospital or organisation e.g. "Trust1".
 - Fixed data type.
 - The LOS_model dataset contains 10 organisations labelled as:
 - * "Trust1"
 - * "Trust2"
 - * "Trust3"

- * "Trust4"
- * "Trust5"
- * "Trust6"
- * "Trust7"
- * "Trust8"
- * "Trust9"
- * "Trust10"
- Allowed values: the data tool only allows selection from the 10 organisations in the LOS_model dataset.
- **Age** : an integer value representing age of the patient.
 - Quantitative data type.
 - Min value: 0
 - Max value: 150
 - Allowed values: only integer values allowed.
 - Example: 45
- **LOS** : length of stay, an integer value representing the days a patient was in hospital.
 - Quantitative data type.
 - Min value: 0
 - Max value: None
 - Allowed values: only integer values allowed.
 - Example: 3
- **Death** : a factor consisting of a string that is either "Alive" or "Died", representing the status of the patient at the end of their hospital stay.
 - Fixed data type.
 - **Alive** : The patient survived the hospital admission.
 - **Died** : The patient died during hospital admission.
 - Allowed values: "Alive" or "Died".
- **consent** : a boolean value representing the consent from the end-user to process and share the data collected with the data capture tool.
 - Fixed data type.
 - **TRUE** : Consent granted by end-user.
 - **FALSE** : Consent not granted by end-user.
 - Allowed values: TRUE or FALSE.

Dataset attributes

Relevant dataset attributes are below:

- **main** : string containing the main description of the dataset.
 - Simple attribute type
- **dictionary** : dataframe object containing the data dictionary.
 - Complex attribute type
 - **variable name** : string vector containing the name of each variable/column.
 - **variable description** : string vector containing the description of each variable/column.
 - **variable options** : string vector containing the data element options for each variable/column.
 - **notes** : string vector containing descriptions of variables and their options.
- **last_edit_date** : date format "YYYY-MM-DD HH:MM:SS" of the last time of edit of the dataset and attributes.
 - Simple attribute type
- **author** : string containing the name of the dataset and dictionary author. E.g. "B213753*"
 - Simple attribute type

Loading R dataset

Loading raw data into R

The following code loads the raw `LOS_model` data from the `NHSRdatasets` package into the R environment and assigns the data to a variable called `LOS`.

```
data(LOS_model)
LOS <- LOS_model
```

Dataset overview

Before getting the data ready for exploratory and statistical analyses, using `glimpse()` and `head()`, the dataset was briefly explored for data types, and for any obvious anomalies.

Glimpse

```
glimpse(LOS)

## Rows: 300
## Columns: 5
## $ ID          <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17~
## $ Organisation <ord> Trust1, Trust2, Trust3, Trust4, Trust5, Trust6, Trust7, T~
## $ Age         <int> 55, 27, 93, 45, 70, 60, 25, 48, 51, 81, 58, 16, 21, 82, 1~
## $ LOS         <int> 2, 1, 12, 3, 11, 7, 4, 4, 7, 1, 4, 3, 1, 9, 12, 1, 4, 3, ~
## $ Death       <int> 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, ~
```

Using `glimpse()`, the dataset can be shown to have 300 rows and 5 columns. The R data classes of the columns can also be seen (see dictionary for more details):

- **ID** - *int* : ID values are of an integer class. This is acceptable as a row identifier and is used later for indexing.
- **Organisation** - *ord* : Organisation values is of an ordered factor class with 10 levels. This is acceptable for the data type it portrays.
- **Age** - *int* : Age values is of an integer class. This is acceptable for the data type it portrays.
- **LOS** - *int* : LOS (length of stay) values is of an integer class. This is acceptable for the data type it portrays.
- **Death** - *int* : Death values is of an integer class. Although statistical model software can work with binary 0/1 int values instead of factor, this variable is recoded into a factor with two levels for increased transparency and easier visualization (see below).

Death Recode

```
LOS <- LOS %>%
  mutate(Death=recode(factor(Death), '1'='Died', '0'='Alive'))
```

The code above recodes the integer values of the raw data into a factor with two levels: “Alive” and “Died”.

Check missing values

```
LOS %>%
  map(is.na) %>%
  map(sum)
```

```
## $ID
## [1] 0
```

```
##
## $Organisation
## [1] 0
##
## $Age
## [1] 0
##
## $LOS
## [1] 0
##
## $Death
## [1] 0
```

There are no counts of missing values (NA).

Check for error values

```
summary(LOS)
```

```
##          ID          Organisation      Age          LOS          Death
##  Min.   : 1.00   Trust1 : 30   Min.    : 5.00   Min.    : 1.000   Alive:247
## 1st Qu.: 75.75   Trust2 : 30   1st Qu.:24.00   1st Qu.: 2.000   Died : 53
##  Median :150.50   Trust3 : 30   Median :54.00   Median : 4.000
##  Mean   :150.50   Trust4 : 30   Mean    :50.66   Mean    : 4.937
## 3rd Qu.:225.25   Trust5 : 30   3rd Qu.:75.25   3rd Qu.: 7.000
##  Max.   :300.00   Trust6 : 30   Max.    :95.00   Max.    :18.000
##                      (Other):120
```

The summary function shows some of the summary statistics of the dataset. Here we can confirm there are no missing values (which would be counted above). In addition, there appears to be no obvious extreme or error values (such as negative values for age or LOS).

Save raw data

As seen from the brief exploration of the dataset above, the data is extremely clean. An output of the final raw dataset is shown below of the first 10 rows of the data.

```
LOS %>%
  head(10) %>%
  kable()
```

ID	Organisation	Age	LOS	Death
1	Trust1	55	2	Alive
2	Trust2	27	1	Alive
3	Trust3	93	12	Alive
4	Trust4	45	3	Died
5	Trust5	70	11	Alive
6	Trust6	60	7	Alive
7	Trust7	25	4	Alive
8	Trust8	48	4	Alive
9	Trust9	51	7	Died
10	Trust10	81	1	Alive

With no further manipulation or transformation, the data is saved as a csv file as ‘/RawData/Los_raw.csv’

```
write_csv(LOS, here("RawData", "LOS_raw.csv"))
```

Splitting data into test and training sets

```
set.seed(777)
testIndex <- createDataPartition(LOS$ID, p = 10/nrow(LOS),
                                list = FALSE,
                                times = 1)
testIndex
```

```
##      Resample1
## [1,]         1
## [2,]        24
## [3,]        60
## [4,]        85
## [5,]        87
## [6,]       137
## [7,]       159
## [8,]       166
## [9,]       175
## [10,]      251
## [11,]      293
## [12,]      300
```

Using the `createDataPartition()` function from the `caret` package, the ID column of the raw dataset is used as the index to split the dataset into test and training sets. The selected ID numbers are shown above to be split into the test dataset.

```
LOS_test <- LOS[testIndex,]
LOS_train <- LOS[-testIndex,]

# Save a row for assessment marking
LOS_test_marker <- LOS_test[1,]
LOS_test <- LOS_test[2:nrow(LOS_test),]
```

The LOS dataset is split using the `testIndex` as shown above. One observation is removed from the test data to be used for assessment marking purposes.

The final test dataset is shown in the table below.

```
LOS_test %>%
  kable()
```

ID	Organisation	Age	LOS	Death
24	Trust4	16	3	Alive
60	Trust10	12	1	Alive
85	Trust5	90	12	Alive
87	Trust7	26	3	Alive
137	Trust7	78	10	Alive
159	Trust9	53	6	Alive
166	Trust6	18	2	Alive
175	Trust5	77	3	Died
251	Trust1	73	1	Died
293	Trust3	81	8	Died
300	Trust10	93	15	Alive

The observation for marking is shown below.

```
LOS_test_marker %>%  
  kable()
```

ID	Organisation	Age	LOS	Death
1	Trust1	55	2	Alive

The training table is shown below (first 10 rows)

```
LOS_train %>%  
  head(10) %>%  
  kable()
```

ID	Organisation	Age	LOS	Death
2	Trust2	27	1	Alive
3	Trust3	93	12	Alive
4	Trust4	45	3	Died
5	Trust5	70	11	Alive
6	Trust6	60	7	Alive
7	Trust7	25	4	Alive
8	Trust8	48	4	Alive
9	Trust9	51	7	Died
10	Trust10	81	1	Alive
11	Trust1	58	4	Alive

The test dataset has one row set aside for marking. The three data sets (test marker, test, training) are then saved in the 'Data' folder.

```
# Save a row for assessment marking  
LOS_test_marker <- LOS_test[1,]  
LOS_test <- LOS_test[2:nrow(LOS_test),]  
  
# Save test and train datasets as csv into 'Data' folder.  
write_csv(LOS_train, here("Data", "LOS_train.csv"))  
write_csv(LOS_test, here("Data", "LOS_test.csv"))  
write_csv(LOS_test_marker, here("Data", "LOS_test_marker.csv"))
```

Data capture tool (python)

Please see the [jupyter notebook](#) to see the **Python source code** and the **functional data capture tool**. The data capture tool was created using Python within a jupyter notebook. The tool shows the source code of data preparation and uses ipython widgets to create interactive widgets for data entry. By using python, each row of data can be added dynamically by using a submit button, and the collected data table is updated automatically. This is a better user experience compared to iterating through cell blocks for each row of data entered. The collected data is saved as a csv file in the 'RawData' folder.

Building Data dictionary in R

```
CollectedData=read_csv(here("RawData", "CollectedDataFinal.csv"))
```

```
## Rows: 12 Columns: 6

## -- Column specification -----
## Delimiter: ","
## chr (2): Organisation, Death
## dbl (3): ID, Age, LOS
## lgl (1): consent

##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.

glimpse(CollectedData)
```

```
## Rows: 12
## Columns: 6
## $ ID          <dbl> 4, 9, 46, 112, 114, 135, 187, 204, 213, 244, 260, 267
## $ Organisation <chr> "Trust3", "Trust9", "Trust6", "Trust2", "Trust4", "Trust5~
## $ Age         <dbl> 0, 51, 49, 19, 23, 93, 33, 27, 48, 76, 11, 79
## $ LOS         <dbl> 0, 7, 7, 3, 3, 10, 4, 2, 5, 9, 1, 10
## $ Death       <chr> "Alive", "Died", "Alive", "Alive", "Alive", "Alive", "Die~
## $ consent     <lgl> TRUE, TRUE, TRUE, TRUE, TRUE, TRUE, TRUE, TRUE, TRU~
```

The Data collected from the collection tool is loaded into R. As shown, there is a consent column which is added during the data collection in the tool.

Create Linker dataframe

```
variable_description <- c(
  'Integer value representing patient number - should be unique and allows refer
  encing rows to the original raw data.',
  'A factor consisting of a string that represents a hospital or organisation
  e.g. "Trust1"',
  'Integer value representing age of the patient.',
  'Integer value representing the days a patient was in hospital.',
  'A factor consisting of a string that is either "Alive" or "Died",
  representing the status of the patient at the end of their hospital stay.',
  'A boolean value representing the consent from the end-user to process and
  share the data collected with the data capture tool.'
)

variable_type <- c(0,1,0,0,1,1)

linker<-build_linker(CollectedData, variable_description, variable_type)
```

A linker dataframe is created using variable_description (string vector describing each column) and variable_type (integer vector, 0 for quantitative and 1 for fixed values). This dataframe will then be used to create a data dictionary

Build dictionary

```
## Enter description for variable 'Age' and option '0 to 93':
## Enter description for variable 'consent' and option 'TRUE':
## Enter description for variable 'Death' and option 'Alive':
## Enter description for variable 'Death' and option 'Died':
## Enter description for variable 'ID' and option '4 to 267':
```



```
## Enter description for variable 'LOS' and option '0 to 10':
## Enter description for variable 'Organisation' and option 'Trust3':
## Enter description for variable 'Organisation' and option 'Trust9':
## Enter description for variable 'Organisation' and option 'Trust6':
## Enter description for variable 'Organisation' and option 'Trust2':
## Enter description for variable 'Organisation' and option 'Trust4':
## Enter description for variable 'Organisation' and option 'Trust5':
## Enter description for variable 'Organisation' and option 'Trust7':
## Enter description for variable 'Organisation' and option 'Trust10':
```

A dictionary is build using `build_dict()`, using the collected data and the linker dataframe. Notes are currently left empty to be manually filled below.

```
# Add notes
dictionary[3,4]<-'Alive: The patient survived the hospital admission.'
dictionary[4,4]<-'Died: The patient died during hospital admission.'
dictionary[7,4]<-'Trust3: Name of hospital/organisation that admitted patient.'
dictionary[8,4]<-'Trust9: Name of hospital/organisation that admitted patient.'
dictionary[9,4]<-'Trust6: Name of hospital/organisation that admitted patient.'
dictionary[10,4]<-'Trust2: Name of hospital/organisation that admitted patient.'
dictionary[11,4]<-'Trust4: Name of hospital/organisation that admitted patient.'
dictionary[12,4]<-'Trust5: Name of hospital/organisation that admitted patient.'
dictionary[13,4]<-'Trust7: Name of hospital/organisation that admitted patient.'
dictionary[14,4]<-'Trust10: Name of hospital/organisation that admitted patient.'
```

Save dictionary

```
write_csv(dictionary, here("RawData", "CollectedData_DataDictionary.csv"))
```

The dictionary saved, and can be re-used for future data collection.

Append dictionary onto dataset

```
# Main string description
main_string <- "This data describes the hospital length of stay (LOS) and outcome
(death/survived) data from the *NHSRdatasets* package collected by the data
capture tool."

complete_CollectedData <- incorporate_attr(my.data = CollectedData,
                                             data.dictionary = dictionary,
                                             main_string = main_string)

# Change the author name
attributes(complete_CollectedData)$author[1]<-"B213753"
```

The data dictionary is now appended onto the attributes of the collected dataset using `incorporate_attr()`, along with a main description and the author name. The attributes can be seen above.

Save completed dataset as R dataset file.

```
save_it(complete_CollectedData, here("RawData", "complete_CollectedData"))
```