

Data_splitting

Hun

2022-03-23

Importing data

Algorithm for splitting variables into training and testing set evenly

```
# My goal is to use quantiles to split data evenly
split <- function(variable, p) {

  #Case for categorical variable or small number of levels of continuous variable
  if (length(unique(variable)) < 10){
    unique <- unique(variable)
    unique_list <- list()

    training_store <- vector()

    testing_store <- vector()

    for (i in 1:length(unique)) {
      unique_list[[i]] <- variable[variable %in% unique[i]]

      #To make sure all variables of training/testing have the same number
      ifelse((i %% 2) == 1 && (length(variable)*p)%1 > 0.5,
        index <- sample(1:length(unique_list[[i]]),
          ceiling(length(unique_list[[i]])*p),replace=FALSE),
        index <- sample(1:length(unique_list[[i]]),
          floor(length(unique_list[[i]])*p),replace=FALSE))
      #To make sure all variables of training/testing have the same number
      ifelse((i %% 2) == 1 && (length(variable)*p)%1 < 0.5,
        index <- sample(1:length(unique_list[[i]]),
          floor(length(unique_list[[i]])*p),replace=FALSE),
        index <- sample(1:length(unique_list[[i]]),
          ceiling(length(unique_list[[i]])*p),replace=FALSE))

      training_store <- c(training_store, unique_list[[i]][index])
      testing_store <- c(testing_store, unique_list[[i]][-index])
    }
    split <- list(training = training_store, testing = testing_store)
  }
}
```

```

else {

  #Case for continuous variable
  variable <- sort(variable)

  #Using quantiles to split data evenly
  smallest <- min(variable)
  first_quantile <- variable[round(length(variable)*0.25, 0)]
  second_quantile <- variable[round(length(variable)*0.5, 0)]
  third_quantile <- variable[round(length(variable)*0.75, 0)]
  largest <- max(variable)

  summary <- c(smallest, first_quantile, second_quantile, third_quantile, largest)

  training_data <- list()
  testing_data <- list()

  training_store <- vector()

  testing_store <- vector()

  for(i in 1:(length(summary)-1)){
    #To make sure all variables of training/testing have the same number
    ifelse(i == (length(summary)-1),
          Q_data <- variable[variable>=summary[i] & variable<=summary[i+1]],
          Q_data <- variable[variable>=summary[i] & variable<summary[i+1]])
    #To make sure all variables of training/testing have the same number
    ifelse((i %% 2) == 1,
          index <- sample(1:length(Q_data), ceiling(length(Q_data)*p), replace=FALSE),
          index <- sample(1:length(Q_data), floor(length(Q_data)*p), replace=FALSE))

    training_data[[i]] <- Q_data[index]

    testing_data[[i]] <- Q_data[-index]

    training_store <- c(training_store, training_data[[i]])

    testing_store <- c(testing_store, testing_data[[i]])
  }
  split <- list(training = training_store, testing = testing_store)
}
return(split)
}

```

```
#Let's check if the algorithm works for continuous variable
```

```
split1 <- split(data$radius_mean, 0.8)
```

```
data_frame(mean = mean(split1$training), sd = sd(split1$training),  
            proportion = length(split1$training)/nrow(data))
```

```
## # A tibble: 1 x 3  
##   mean    sd proportion  
##   <dbl> <dbl>     <dbl>  
## 1  14.2  3.56      0.800
```

```
data_frame(mean = mean(split1$testing), sd = sd(split1$testing),  
            proportion = length(split1$testing)/nrow(data))
```

```
## # A tibble: 1 x 3  
##   mean    sd proportion  
##   <dbl> <dbl>     <dbl>  
## 1  14.0  3.40      0.200
```

```
#Looks good
```

```
#Let's check if the algorithm works for categorical variable
```

```
split2 <- split(data$diagnosis, 0.8)
```

```
data_frame(mean = mean(split2$training), sd = sd(split2$training),  
            proportion = length(split2$training)/length(data$radius_mean))
```

```
## # A tibble: 1 x 3  
##   mean    sd proportion  
##   <dbl> <dbl>     <dbl>  
## 1  0.371 0.484      0.800
```

```
data_frame(mean = mean(split2$testing), sd = sd(split2$testing),  
            proportion = length(split2$testing)/length(data$radius_mean))
```

```
## # A tibble: 1 x 3  
##   mean    sd proportion  
##   <dbl> <dbl>     <dbl>  
## 1  0.377 0.487      0.200
```

```
#Looks good
```

Applying algorithm to get the training/testing data frame of entire data

```
data_split <- function(data,split, p){  
  data_split <- map(data, split, p)  
  training_list <- list()
```

```

testing_list <- list()
for (i in 1:length(data_split)) {

  training_list[[i]] <- data_split[[i]]$training
  testing_list[[i]] <- data_split[[i]]$testing
}

names(training_list) <- names(data_split)
names(testing_list) <- names(data_split)

training <- dplyr::bind_rows(training_list) %>% data.frame()
testing <- dplyr::bind_rows(testing_list) %>% data.frame()

return(list(training, testing))
}

```

Combining result together to make it reader-frindly

```

trainging_result <-
  skimr::skim_without_charts(data_split(data,split, 0.8)[1]) %>% data.frame() %>%
    select(2,5,6) %>%
    rename(training_mean = numeric.mean, training_sd = numeric.sd) %>%
    mutate_if(is.numeric, ~round(.x, digits = 3)) %>%
    mutate_if(is.numeric, ~format(.x, scientific = FALSE))

testing_result <-
  skimr::skim_without_charts(data_split(data,split, 0.8)[2]) %>%
    select(2,5,6) %>%
    rename(testing_mean = numeric.mean, testing_sd = numeric.sd) %>%
    mutate_if(is.numeric, ~round(.x, digits = 3)) %>%
    mutate_if(is.numeric, ~format(.x, scientific = FALSE))

trainging_result %>% left_join(testing_result, by = "skim_variable") %>% kable()

```

skim_variable	training_mean	training_sd	testing_mean	testing_sd
id	27573499.703	112500444.331	20935774.711	86027240.237
diagnosis	0.371	0.484	0.377	0.487
radius_mean	14.111	3.484	13.992	3.278
texture_mean	19.261	4.177	19.102	4.410
smoothness_mean	0.097	0.014	0.096	0.015
compactness_mean	0.104	0.052	0.103	0.050
concave_points_mean	0.049	0.039	0.048	0.036
symmetry_mean	0.181	0.027	0.181	0.029
fractal_dimension_mean	0.063	0.007	0.063	0.008
radius_se	0.399	0.260	0.401	0.239
texture_se	1.212	0.545	1.230	0.551
smoothness_se	0.007	0.003	0.007	0.004
compactness_se	0.026	0.018	0.024	0.015
concavity_se	0.032	0.032	0.030	0.021
concave_points_se	0.012	0.006	0.011	0.005
symmetry_se	0.021	0.008	0.021	0.009
fractal_dimension_se	0.004	0.003	0.004	0.002
radius_worst	16.251	4.844	16.306	4.711
smoothness_worst	0.132	0.023	0.133	0.023
compactness_worst	0.254	0.156	0.251	0.158
concavity_worst	0.273	0.209	0.261	0.200
symmetry_worst	0.290	0.061	0.287	0.055
fractal_dimension_worst	0.084	0.018	0.083	0.016

training_nrow	testing_nrow	eighty_percent_data_nrow	twenty_percent_data_nrow
455	114	455.2	113.8