

SHETH L.U.J AND SIR M.V. COLLEGE

SUBJECT :- DATA ANALYSIS WITH SAS/SPSS/R

PRACTICAL – 9

AIM:- Performing text manipulation using str_sub(), str_split() (R).
import dataset.

OUTPUT:-

```
R - R 4.5.2 - ~/r
> library(stringr)
> library(tidyverse)
> library(dplyr)
> library(readxl)
>
> # 1. IMPORT DATASET
> spotify <- read_excel("C:\\\\Users\\\\info\\\\Downloads\\\\spotify_data_clean.xlsx")
>
> print("--- original Dataset (First 5 Rows) ---")
[1] "--- original Dataset (First 5 Rows) ---"
> print(head(spotify, 5))
# A tibble: 5 x 4
  track_id track_name track_number track_popularity
  <chr>     <chr>           <dbl>            <dbl>
1 3EJ5Lyek0imTrippy Ma...     4             0 TRUE
2 1oQw6G22iWmU.OMG!       1             0 TRUE
3 7ndkjzoiYf1...        1             4 TRUE
4 67rw02l7083q...        8             30 TRUE
5 1XpxTfRrjs...       2             0 TRUE
# i more variables: album_name <chr>, album_release_date <dttm>, album_total_tracks <dbl>, album_type <chr>,
# track_duration_min <dbl>
>
> # 2. USING str_sub() (Substring)
> spotify$ID_code <- str_sub(spotify$track_id, 1, 6)
> spotify$Year <- str_sub(as.character(spotify$album_release_date), -4, -1)
>
> print("--- Data after str_sub() ---")
[1] "... Data after str_sub() ---"
> print(spotify %>% select(track_id, ID_Code, album_release_date, Year) %>% head(5))
# A tibble: 5 x 4
  track_id ID_Code album_release_date Year
  <chr>    <chr>    <dttm>          <chr>
1 3EJ5Lyek0imTrippy Ma... 3EJ5Ly 2023-10-31 00:00:00 0-31
2 1oQw6G22iWmU.OMG!      1oQw6G 2023-10-31 00:00:00 0-31
3 7ndkjzoiYf1...        7ndkjz 2023-10-31 00:00:00 0-31
4 67rw02l7083q...        67rw02 2023-10-31 00:00:00 0-31
5 1XpxTfRrjs...        1XpxTf 2023-10-30 00:00:00 0-30
>
> # 3. USING str_split() (Split string)
> # Method A: Basic split (list output)
> genre_list <- str_split(spotify$artist_genres, ",")
> print("--- Basic Split Output (List format) ---")
[1] "... Basic Split Output (List format) ---"
> print(genre_list[[1]]) # show first artist's genre list
[1] "moombahton"
>
> # Method B: Split Fixed (Matrix form)
> library(stringr)
> library(tidyverse)
> library(dplyr)
> library(readxl)
>
> # 1. IMPORT DATASET
> spotify <- read_excel("C:\\\\Users\\\\info\\\\Downloads\\\\spotify_data_clean.xlsx")
>
> print("--- original Dataset (First 5 Rows) ---")
[1] "--- original Dataset (First 5 Rows) ---"
> print(head(spotify, 5))
# A tibble: 5 x 4
  track_id track_name track_number track_popularity
  <chr>     <chr>           <dbl>            <dbl>
1 3EJ5Lyek0imTrippy Ma...     4             0 TRUE
2 1oQw6G22iWmU.OMG!       1             0 TRUE
3 7ndkjzoiYf1...        1             4 TRUE
4 67rw02l7083q...        8             30 TRUE
5 1XpxTfRrjs...       2             0 TRUE
# i more variables: album_name <chr>, album_release_date <dttm>, album_total_tracks <dbl>, album_type <chr>,
# track_duration_min <dbl>
>
> # 2. USING str_sub() (Substring)
> spotify$ID_code <- str_sub(spotify$track_id, 1, 6)
> spotify$Year <- str_sub(as.character(spotify$album_release_date), -4, -1)
>
> print("--- Data after str_sub() ---")
[1] "... Data after str_sub() ---"
> print(spotify %>% select(track_id, ID_Code, album_release_date, Year) %>% head(5))
# A tibble: 5 x 4
  track_id ID_Code album_release_date Year
  <chr>    <chr>    <dttm>          <chr>
1 3EJ5Lyek0imTrippy Ma... 3EJ5Ly 2023-10-31 00:00:00 0-31
2 1oQw6G22iWmU.OMG!      1oQw6G 2023-10-31 00:00:00 0-31
3 7ndkjzoiYf1...        7ndkjz 2023-10-31 00:00:00 0-31
4 67rw02l7083q...        67rw02 2023-10-31 00:00:00 0-31
5 1XpxTfRrjs...        1XpxTf 2023-10-30 00:00:00 0-30
```

```
R - R 4.5.2 - ~/r
>
> # 3. USING str_split() (Split string)
> # Method A: Basic split (list output)
> genre_list <- str_split(spotify$artist_genres, ",")
> print("--- Basic Split Output (List format) ---")
[1] "... Basic Split Output (List format) ---"
> print(genre_list[[1]]) # show first artist's genre list
[1] "moombahton"
>
> # Method B: Split Fixed (Matrix form)
> library(stringr)
> library(tidyverse)
> library(dplyr)
> library(readxl)
>
> # 1. IMPORT DATASET
> spotify <- read_excel("C:\\\\Users\\\\info\\\\Downloads\\\\spotify_data_clean.xlsx")
>
> print("--- original Dataset (First 5 Rows) ---")
[1] "--- original Dataset (First 5 Rows) ---"
> print(head(spotify, 5))
# A tibble: 5 x 4
  track_id track_name track_number track_popularity
  <chr>     <chr>           <dbl>            <dbl>
1 3EJ5Lyek0imTrippy Ma...     4             0 TRUE
2 1oQw6G22iWmU.OMG!       1             0 TRUE
3 7ndkjzoiYf1...        1             4 TRUE
4 67rw02l7083q...        8             30 TRUE
5 1XpxTfRrjs...       2             0 TRUE
# i more variables: album_name <chr>, album_release_date <dttm>, album_total_tracks <dbl>, album_type <chr>,
# track_duration_min <dbl>
>
> # 2. USING str_sub() (Substring)
> spotify$ID_code <- str_sub(spotify$track_id, 1, 6)
> spotify$Year <- str_sub(as.character(spotify$album_release_date), -4, -1)
>
> print("--- Data after str_sub() ---")
[1] "... Data after str_sub() ---"
> print(spotify %>% select(track_id, ID_Code, album_release_date, Year) %>% head(5))
# A tibble: 5 x 4
  track_id ID_Code album_release_date Year
  <chr>    <chr>    <dttm>          <chr>
1 3EJ5Lyek0imTrippy Ma... 3EJ5Ly 2023-10-31 00:00:00 0-31
2 1oQw6G22iWmU.OMG!      1oQw6G 2023-10-31 00:00:00 0-31
3 7ndkjzoiYf1...        7ndkjz 2023-10-31 00:00:00 0-31
4 67rw02l7083q...        67rw02 2023-10-31 00:00:00 0-31
5 1XpxTfRrjs...        1XpxTf 2023-10-30 00:00:00 0-30
```

SHETH L.U.J AND SIR M.V. COLLEGE

SUBJECT :- DATA ANALYSIS WITH SAS/SPSS/R

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Source Terminal Background Jobs
[R - R 4.5.2 - ~]
> # 2. USING str_sub() (substring)
> spotify$Id_code <- str_sub(spotify$track_id, 1, 6)
> spotify$year <- str_sub(as.character(spotify$album_release_date), -4, -1)

> print("---- Data after str_sub() ----")
[1] "---- Data after str_sub() ----"
> print(spotify %>% select(track_id, Id_Code, album_release_date, year) %>% head(5))
# A tibble: 5 x 4
  track_id      Id_Code album_release_date   year
<chr>        <chr>    <dttm>       <dbl>
1 3EJ5LyeK0mtf5rBmzI 3EJ5L 2023-10-31 00:00:00 0.31
2 3n0u6C2iwuuiq1pp270o 1en66C 2023-10-31 00:00:00 0.31
3 7ndkjzoitYfirx9ct8pmu 7ndkjz 2023-10-31 00:00:00 0.31
4 67rw02l7or83qpo5vw5w 67rw02 2023-10-31 00:00:00 0.31
5 1xpttfRBrjsppwOinUu2jf 1xpttf 2023-10-30 00:00:00 0.30

> # 3. USING str_split() (Split String)
> # Method A: Basic split output
> genre_list <- str_split(spotify$artist_genres, ", ")
> print("---- Basic split output (List format) ----")
[1] "---- Basic split output (List format) ----"
> print(genre_list[[1]]) # show first artist's genre list
[1] "moombahton"
>
> # Method B: Split Fixed (Matrix form)
> genre_matrix <- str_split(spotify$artist_genres, ",", simplify = TRUE)
>
> spotify$Genre_1 <- genre_matrix[, 1]
> spotify$Genre_2 <- genre_matrix[, 2]

> print("---- Data after str_split() (Manual Assignment) ----")
[1] "---- Data after str_split() (Manual Assignment) ----"
> print(spotify %>% select(artist_genres, Genre_1, Genre_2) %>% head(5))
# A tibble: 5 x 3
  artist_genres      Genre_1      Genre_2
<chr>            <chr>        <chr>
1 moombahton        moombahton    
2 country hip hop, southern hip hop country hip hop " southern hip hop"
3 N/A                N/A          ...
4 moombahton        moombahton    
5 dark r&b         dark r&b     ...

> # 4. BONUS: Using 'separate' to split track_name
> tidy_spotify %>% separate(track_name, into = c("title", "info"), sep = " - ", fill = "right")
Warning message:
Expected 2 pieces. Additional pieces discarded in 11 rows [581, 2003, 3739, 3983, 3984, 4261, 4570, 5668, 6673, 7122, 8548].
>
> print("---- Bonus: The 'separate' function (Track Title split) ----")
[1] "---- Bonus: The 'separate' function (Track Title split) ----"
> print(tidy_spotify %>% select(title, info) %>% head(5))
# A tibble: 5 x 2
  title           info
<chr>          <chr>
1 Party Mane (ft. Project Pat) NA
2 OMG!           NA
3 Hard 2 Find    NA
4 Still Get Like That (ft. Project Pat & Starrah) NA
5 ride me like a harley  NA

```

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Source Terminal Background Jobs
[R - R 4.5.2 - ~]
> print("---- Basic Split Output (List format) ----")
[1] "---- Basic Split Output (List format) ----"
> print(genre_list[[1]]) # show first artist's genre list
[1] "moombahton"
>
> # Method B: Split Fixed (Matrix form)
> genre_matrix <- str_split(spotify$artist_genres, ",", simplify = TRUE)
>
> spotify$Genre_1 <- genre_matrix[, 1]
> spotify$Genre_2 <- genre_matrix[, 2]

> print("---- Data after str_split() (Manual Assignment) ----")
[1] "---- Data after str_split() (Manual Assignment) ----"
> print(spotify %>% select(artist_genres, Genre_1, Genre_2) %>% head(5))
# A tibble: 5 x 3
  artist_genres      Genre_1      Genre_2
<chr>            <chr>        <chr>
1 moombahton        moombahton    
2 country hip hop, southern hip hop country hip hop " southern hip hop"
3 N/A                N/A          ...
4 moombahton        moombahton    
5 dark r&b         dark r&b     ...

> # 4. BONUS: Using 'separate' to split track_name
> tidy_spotify %>% separate(track_name, into = c("title", "info"), sep = " - ", fill = "right")
Warning message:
Expected 2 pieces. Additional pieces discarded in 11 rows [581, 2003, 3739, 3983, 3984, 4261, 4570, 5668, 6673, 7122, 8548].
>
> print("---- Bonus: The 'separate' function (Track Title split) ----")
[1] "---- Bonus: The 'separate' function (Track Title split) ----"
> print(tidy_spotify %>% select(title, info) %>% head(5))
# A tibble: 5 x 2
  title           info
<chr>          <chr>
1 Party Mane (ft. Project Pat) NA
2 OMG!           NA
3 Hard 2 Find    NA
4 Still Get Like That (ft. Project Pat & Starrah) NA
5 ride me like a harley  NA

```

SHETH L.U.J AND SIR M.V. COLLEGE
SUBJECT :- DATA ANALYSIS WITH SAS/SPSS/R

SUMEET - S124