

PRACTICAL NO 9

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

CODE:

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)
S090_PRACTICAL6.R* S090_PRACTICAL7.R* S090_PRACTICAL8.R* S090_PRACTICAL9.R* flipkart_commerce_sample
1 #9. Performing text manipulation using str_sub(), str_split() (R). import dataset.
2 # =====
3 # R script: Text Manipulation with stringr
4 # Functions: str_sub(), str_split()
5 # =====
6
7 # Load necessary library
8 install.packages("stringr")
9 install.packages("tidyr") # for separating columns after splitting
10 library(stringr)
11 library(tidyr)
12 library(dplyr)
13
14 # 1. CREATE DATASET
15
16 # We are adding a 'SKU' column which is perfect for text manipulation practice.
17 # Format: "CATEGORY-PRODUCTID-YEAR" (e.g., "ELEC-5548-2023")
18
19
20 retail_data <- data.frame(
21   SKU = c("ELEC-5548-2023", "HOME-3045-2022", "CLOT-4004-2023", "ELEC-4808-2021", "HOME-1817-2023"),
22   Description = c("Electronics - Smart TV", "Home - Blender", "Clothing - TShirt", "Electronics - Laptop", "Home - Sofa"),
23   Price = c(500, 45, 20, 900, 300)
24 )
25
26 print("--- Original Dataset ---")
27 print(retail_data)
28
29 #
30 # 2. USING str_sub() (Substring)
31 # =====
32 # Scenario: We want to extract specific parts of the SKU based on position.
33 # Syntax: str_sub(string, start, end)
34
79:1 (Untitled) R Script

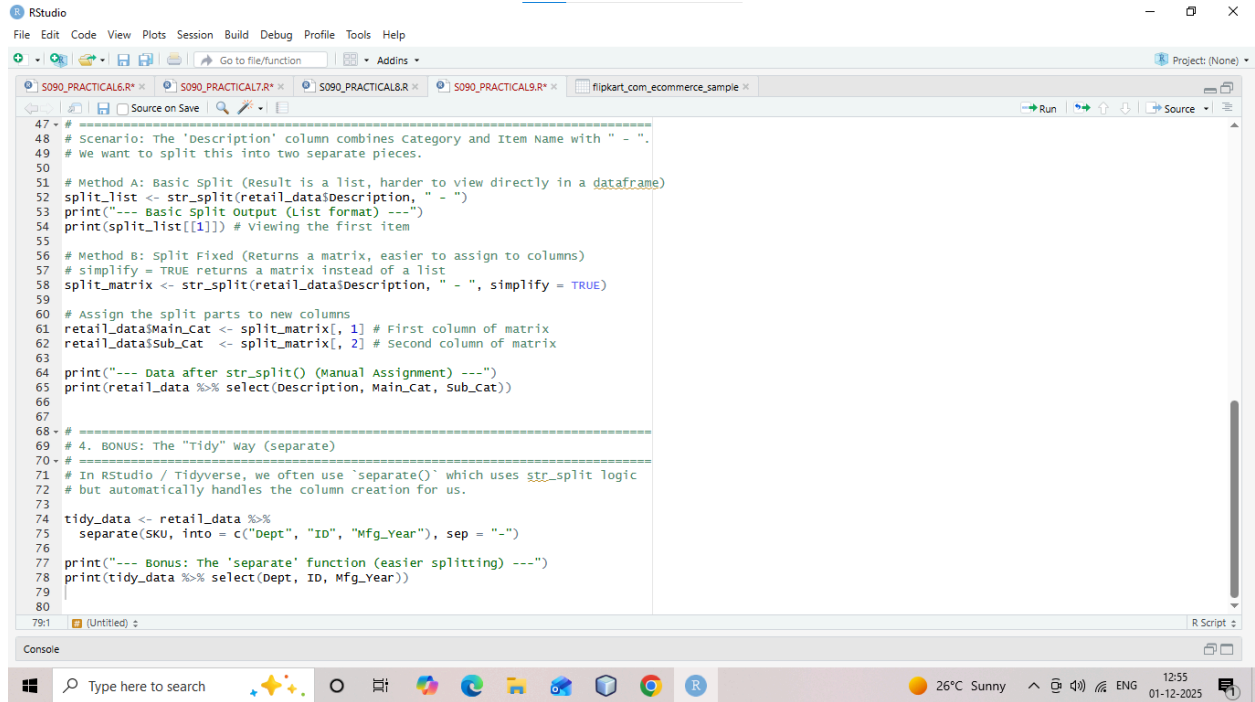
```

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)
S090_PRACTICAL6.R* S090_PRACTICAL7.R* S090_PRACTICAL8.R* S090_PRACTICAL9.R* flipkart_commerce_sample
35 # Example A: Extract the first 4 characters to get the Category code
36 retail_data$Category_Code <- str_sub(retail_data$SKU, 1, 4)
37
38 # Example B: Extract the last 4 characters to get the Year
39 # We can use negative numbers to count from the end of the string.
40 retail_data$Year <- str_sub(retail_data$SKU, -4, -1)
41
42 print("--- Data after str_sub() ---")
43 print(retail_data %>% select(SKU, Category_Code, Year))
44
45 #
46 # 3. USING str_split() (Split string)
47 # =====
48 # Scenario: The 'Description' column combines Category and Item Name with " - ".
49 # We want to split this into two separate pieces.
50
51 # Method A: Basic Split (Result is a list, harder to view directly in a dataframe)
52 split_list <- str_split(retail_data$Description, " - ")
53 print("--- Basic split output (List format) ---")
54 print(split_list[[1]]) # Viewing the first item
55
56 # Method B: Split Fixed (Returns a matrix, easier to assign to columns)
57 # simplify = TRUE returns a matrix instead of a list
58 split_matrix <- str_split(retail_data$Description, " - ", simplify = TRUE)
59
60 # Assign the split parts to new columns
61 retail_data$Main_cat <- split_matrix[, 1] # First column of matrix
62 retail_data$Sub_cat <- split_matrix[, 2] # Second column of matrix
63
64 print("--- Data after str_split() (Manual Assignment) ---")
65 print(retail_data %>% select(Description, Main_cat, Sub_cat))
66
67
68 #
79:1 (Untitled) R Script

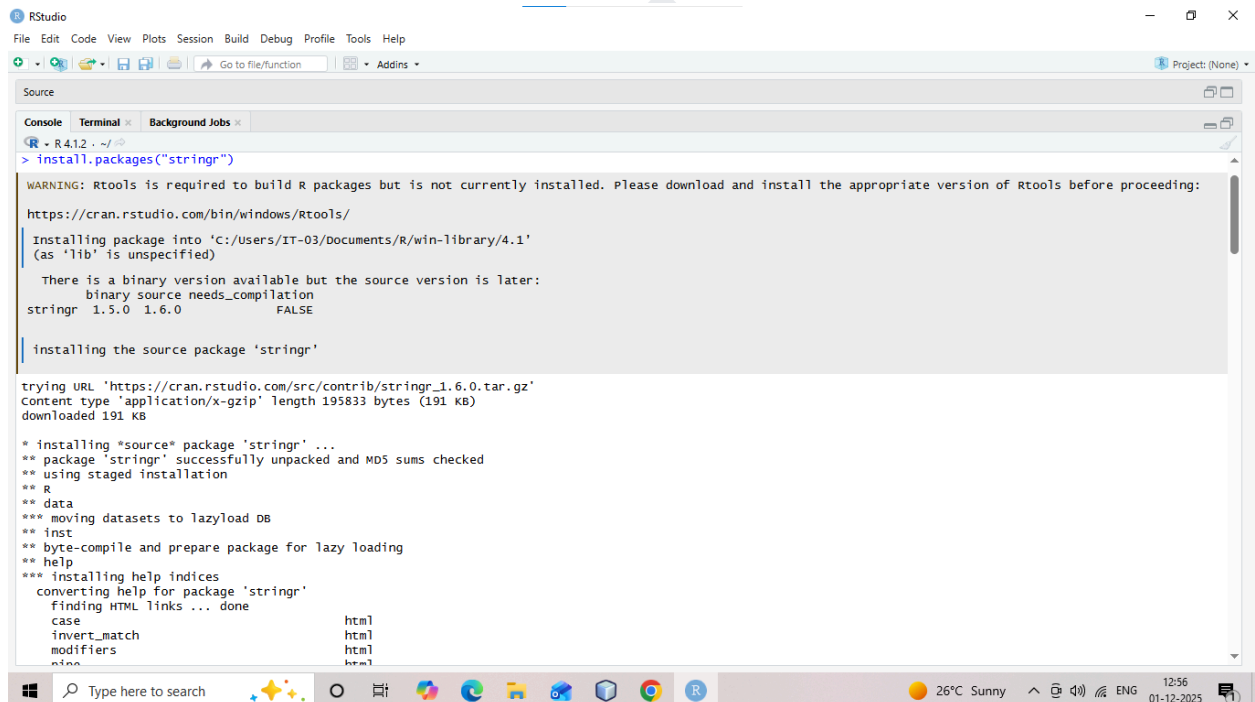
```

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



```
47 # =====
48 # Scenario: The 'description' column combines category and Item Name with " - ".
49 # We want to split this into two separate pieces.
50
51 # Method A: Basic Split (Result is a list, harder to view directly in a dataframe)
52 split_list <- str_split(retail_data$description, " - ")
53 print("--- Basic Split output (List format) ---")
54 print(split_list[[1]]) # Viewing the first item
55
56 # Method B: Split Fixed (Returns a matrix, easier to assign to columns)
57 # simplify = TRUE returns a matrix instead of a list
58 split_matrix <- str_split(retail_data$description, " - ", simplify = TRUE)
59
60 # Assign the split parts to new columns
61 retail_data$Main_Cat <- split_matrix[, 1] # First column of matrix
62 retail_data$Sub_Cat <- split_matrix[, 2] # Second column of matrix
63
64 print("--- Data after str_split() (Manual Assignment) ---")
65 print(retail_data %>% select(Description, Main_Cat, Sub_Cat))
66
67 # =====
68 # 4. BONUS: The "Tidy" Way (separate)
69 # =====
70 # In RStudio / Tidyverse, we often use 'separate()' which uses str_split logic
71 # but automatically handles the column creation for us.
72
73 tidy_data <- retail_data %>%
74   separate(SKU, into = c("Dept", "ID", "Mfg_Year"), sep = "-")
75
76 print("--- Bonus: The 'separate' function (easier splitting) ---")
77 print(tidy_data %>% select(Dept, ID, Mfg_Year))
78
79
80
```

OUTPUT:



```
> install.packages("stringr")

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:
https://cran.rstudio.com/bin/windows/Rtools/

Installing package into 'C:/Users/IT-03/Documents/R/win-library/4.1'
(as 'lib' is unspecified)

There is a binary version available but the source version is later:
  binary source needs_compilation
stringr  1.5.0 1.6.0          FALSE

Installing the source package 'stringr'

trying URL 'https://cran.rstudio.com/src/contrib/stringr_1.6.0.tar.gz'
Content type 'application/x-gzip' length 195833 bytes (191 KB)
downloaded 191 KB

* installing *source* package 'stringr' ...
** package 'stringr' successfully unpacked and MD5 sums checked
** using staged installation
** R
** data
*** moving datasets to lazyload DB
** inst
** byte-compile and prepare package for lazy loading
** help
*** installing help indices
converting help for package 'stringr'
  finding HTML links ... done
   case
invert_match
modifiers
  html
  html
  html
  html
```

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

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Terminal Background Jobs

R - R 4.1.2 - ~/...
** testing if installed package can be loaded from final location
*** arch - i386
*** arch - x64
** testing if installed package keeps a record of temporary installation path
* DONE (stringr)

The downloaded source packages are in
'c:\Users\IT-03\AppData\Local\Temp\RtmpOE9ZLb\downloaded_packages'
> install.packages("tidyr") # for separating columns after splitting

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:
https://cran.rstudio.com/bin/windows/Rtools/

Installing package into 'c:\Users\IT-03\Documents\R\win-library\4.1'
(as 'lib' is unspecified)
also installing the dependency 'rlang'

There are binary versions available but the source versions are later:
  binary source needs_compilation
rlang  1.1.0 1.1.6 TRUE
tidyr  1.3.0 1.3.1 TRUE

Binaries will be installed

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.1/rlang_1.1.0.zip'
content type 'application/zip' length 1710397 bytes (1.6 MB)
downloaded 1.6 MB

trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.1/tidyr_1.3.0.zip'
content type 'application/zip' length 1440051 bytes (1.4 MB)
downloaded 1.4 MB

package 'rlang' successfully unpacked and MD5 sums checked
package 'tidyr' successfully unpacked and MD5 sums checked
```

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Terminal Background Jobs

R - R 4.1.2 - ~/...

package 'rlang' successfully unpacked and MD5 sums checked
package 'tidyr' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
c:\Users\IT-03\AppData\Local\Temp\RtmpOE9ZLb\downloaded_packages
> library(stringr)
> library(tidyr)

warning message:
package 'tidyr' was built under R version 4.1.3

> library(dplyr)

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':
  filter, lag

The following objects are masked from 'package:base':
  intersect, setdiff, setequal, union

warning message:
package 'dplyr' was built under R version 4.1.3

> retail_data <- data.frame(
+   SKU = c("ELEC-5548-2023", "HOME-3045-2022", "CLOT-4004-2023", "ELEC-4808-2021", "HOME-1817-2023"),
+   Description = c("Electronics - Smart TV", "Home - Blender", "Clothing - Tshirt", "Electronics - Laptop", "Home - Sofa"),
+   Price = c(500, 45, 20, 900, 300)
+ )
> print("--- Original Dataset ---")
[1] "--- Original Dataset ---"
```

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

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Terminal Background Jobs

R 4.1.2 ~\...
1 ELEC-5548-2023 ELEC 2023
2 HOME-3045-2022 HOME 2022
3 CLOT-4004-2023 CLOT 2023
4 ELEC-4808-2021 ELEC 2021
5 HOME-1817-2023 HOME 2023

> # Method A: Basic split (Result is a list, harder to view directly in a dataframe)
> split_list <- str_split(retail_data$Description, " - ")
> print("---- Basic split output (List format) ----")
[1] "---- Basic split output (List format) ----"
> print(split_list[[1]]) # viewing the first item
[1] "Electronics" "Smart TV"
> # Method B: Split Fixed (Returns a matrix, easier to assign to columns)
> # simplify = TRUE returns a matrix instead of a list
> split_matrix <- str_split(retail_data$Description, " - ", simplify = TRUE)
> # Assign the split parts to new columns
> retail_data$Main_Cat <- split_matrix[, 1] # First column of matrix
> retail_data$Sub_Cat <- split_matrix[, 2] # Second column of matrix
> print("---- Data after str_split() (Manual Assignment) ----")
[1] "---- Data after str_split() (Manual Assignment) ----"
> print(retail_data %>% select(Description, Main_Cat, Sub_Cat))
  Description Main_Cat Sub_Cat
1 Electronics - Smart TV Electronics Smart TV
2 Home - Blender Home Blender
3 Clothing - Tshirt Clothing Tshirt
4 Electronics - Laptop Electronics Laptop
5 Home - Sofa Home Sofa

> tidy_data <- retail_data %>%
+ separate(SKU, into = c("Dept", "ID", "Mfg_Year"), sep = "-")
> print("---- Bonus: The 'separate' function (easier splitting) ----")
[1] "---- Bonus: The 'separate' function (easier splitting) ----"
> print(tidy_data %>% select(Dept, ID, Mfg_Year))
  Dept ID Mfg_Year
1 ELEC 5548 2023
2 HOME 3045 2022
3 CLOT 4004 2023
4 ELEC 4808 2021
```