# Sheth L.U.J & Sir M.V College
## SAS/SPSS/R Programming
## Practical No. 11 to 15

Aim: 11. Reshaping data using pivot_longer() and pivot_wider() (R).





Jeetesh Shirganoor

S125

Aim: 12. Combining datasets vertically (concatenation) using rbind() (R).

Aim: 13 Identifying and handling duplicates using distinct() (R).





Jeetesh Shirganoor
S125

# Sheth L.U.J & Sir M.V College
## SAS/SPSS/R Programming
## Practical No. 11 to 15

Aim: 14 Extracting date components using lubridate:: functions (R).



```r
> #14 Extracting date components using lubridate:: functions (R).
> # 1. Load Libraries
> library(lubridate)
> library(dplyr)
>
> # 2. Create Sample Data (Date as text)
> dates_df <- data.frame(
+   Event_ID = 1:4,
+   Date_String = c("2023-01-15", "2023-10-31", "2024-02-29", "2024-12-25")
+ )
>
> # 3. Parse and Extract Components
> processed_data <- dates_df %>%
+   mutate(
+     Actual_Date = ymd(Date_String),               # convert text to date
+     Year_Num = year(Actual_Date),                 # extract year
+     Month_Num = month(Actual_Date),               # extract month number
+     Month_Name = month(Actual_Date, label = TRUE),  # month short name
+     Day_Num = day(Actual_Date),                   # day of month
+     Weekday_Num = wday(Actual_Date),              # weekday number
+     Weekday_Name = wday(Actual_Date, label = TRUE, abbr = FALSE), # weekday name
+     Quarter = quarter(Actual_Date),               # quarter number
+     Day_of_Year = yday(Actual_Date)               # day number in year
+   )
>
> print("--- Extracted Date Components ---")
[1] "--- Extracted Date Components ---"
> print(processed_data)
  Event_ID Date_String Actual_Date Year_Num Month_Num Month_Name Day_Num Weekday_Num Weekday_Name Quarter Day_of_Year
1        1  2023-01-15  2023-01-15     2023         1        Jan      15           1       Sunday       1          15
2        2  2023-10-31  2023-10-31     2023        10        Oct      31           3      Tuesday       4         304
3        3  2024-02-29  2024-02-29     2024         2        Feb      29           5     Thursday       1          60
4        4  2024-12-25  2024-12-25     2024        12        Dec      25           4    Wednesday       4         360
>
> # 4. System Date-Time (Now)
> current_time <- now()
> print("--- Current System Time ---")
[1] "--- Current System Time ---"
> print(paste("Year:", year(current_time)))
[1] "Year: 2025"
> print(paste("Hour:" hour(current time)))
```



```r
> # 2. Create Sample Data (Date as text)
> dates_df <- data.frame(
+   Event_ID = 1:4,
+   Date_String = c("2023-01-15", "2023-10-31", "2024-02-29", "2024-12-25")
+ )
>
> # 3. Parse and Extract Components
> processed_data <- dates_df %>%
+   mutate(
+     Actual_Date = ymd(Date_String),               # convert text to date
+     Year_Num = year(Actual_Date),                 # extract year
+     Month_Num = month(Actual_Date),               # extract month number
+     Month_Name = month(Actual_Date, label = TRUE),  # month short name
+     Day_Num = day(Actual_Date),                   # day of month
+     Weekday_Num = wday(Actual_Date),              # weekday number
+     Weekday_Name = wday(Actual_Date, label = TRUE, abbr = FALSE), # weekday name
+     Quarter = quarter(Actual_Date),               # quarter number
+     Day_of_Year = yday(Actual_Date)               # day number in year
+   )
>
> print("--- Extracted Date Components ---")
[1] "--- Extracted Date Components ---"
> print(processed_data)
  Event_ID Date_String Actual_Date Year_Num Month_Num Month_Name Day_Num Weekday_Num Weekday_Name Quarter Day_of_Year
1        1  2023-01-15  2023-01-15     2023         1        Jan      15           1       Sunday       1          15
2        2  2023-10-31  2023-10-31     2023        10        Oct      31           3      Tuesday       4         304
3        3  2024-02-29  2024-02-29     2024         2        Feb      29           5     Thursday       1          60
4        4  2024-12-25  2024-12-25     2024        12        Dec      25           4    Wednesday       4         360
>
> # 4. System Date-Time (Now)
> current_time <- now()
> print("--- Current System Time ---")
[1] "--- Current System Time ---"
> print(paste("Year:", year(current_time)))
[1] "Year: 2025"
> print(paste("Hour:", hour(current_time)))
[1] "Hour: 11"
> print(paste("Minute:", minute(current_time)))
[1] "Minute: 25"
>
```

Jeetesh Shirganoor

S125

# Sheth L.U.J & Sir M.V College
## SAS/SPSS/R Programming
## Practical No. 11 to 15

Aim: 15 Generating basic summaries using str() or summary() (R).





Jeetesh Shirganoor

S125