# MVLU COLLEGE
## Data analysis with SAS/SSPR/R
## PRACTICAL NO.13

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

## INPUT:

```
library(dplyr)

amazon_df <- read.csv("D:/S119/DATA ANAN/Amazon.csv", na.strings = c("", "NA"))

print("--- 1. Original Dataset (Note rows) ---")
print(amazon_df)

duplicates_report <- amazon_df %>%
  group_by(OrderID, CustomerName, ProductName) %>%
  count() %>%
  filter(n > 1)

print("--- 2. Identification Report (Rows that are duplicated) ---")
print(duplicates_report)

clean_exact <- amazon_df %>%
  distinct()

print("--- 3. Removed Exact Duplicates (distinct) ---")
print(clean_exact)

unique_customers <- amazon_df %>%
  distinct(CustomerName, .keep_all = TRUE)

print("--- 4. Unique Customers Only (Partial Duplicates removed) ---")
print(unique_customers)
```
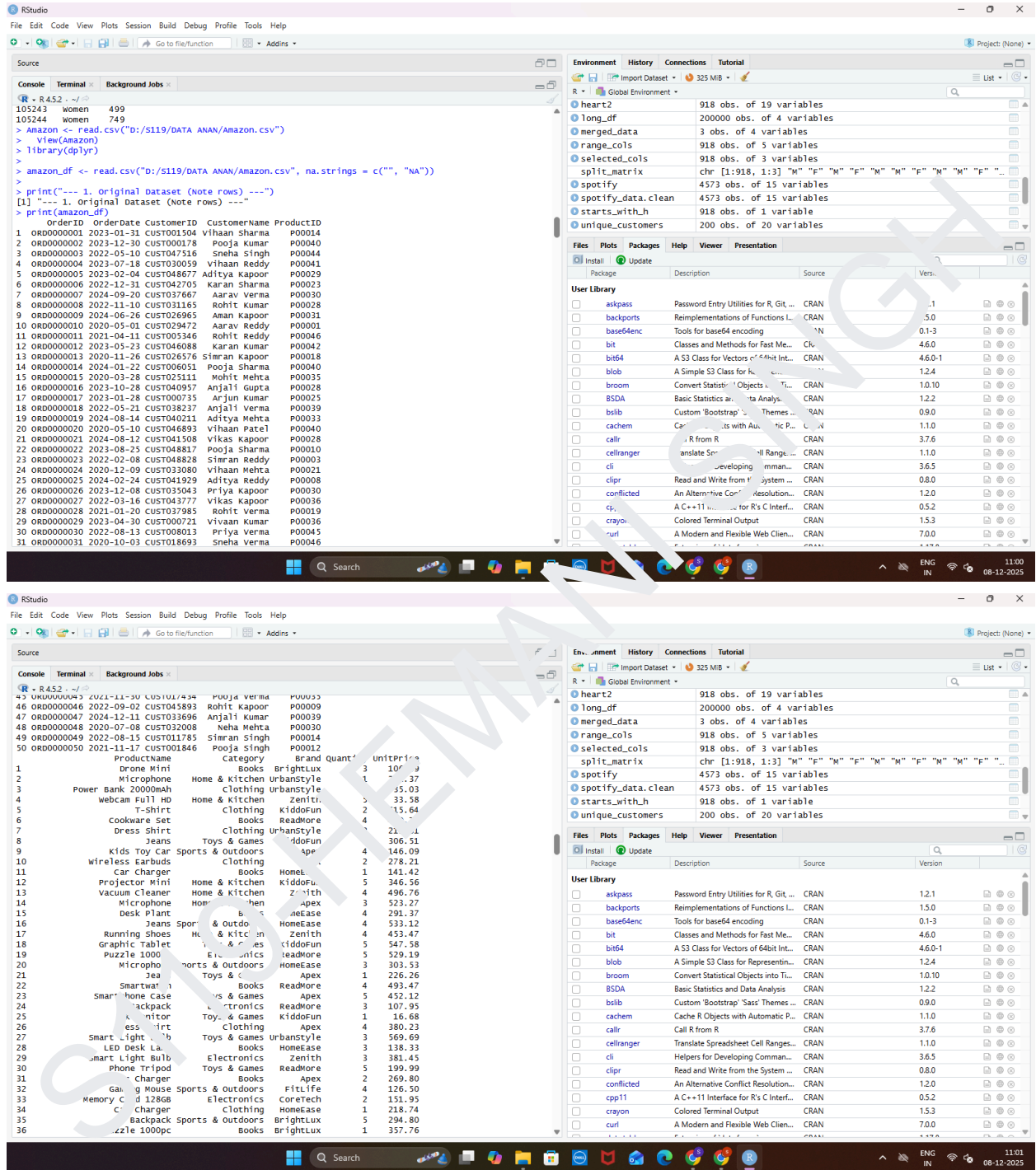
## OUTPUT:

**NAME:HEMANI SINGH**
**ROLL NO:S119**

# MVLU COLLEGE
# Data analysis with SAS/SSPR/R
# PRACTICAL NO.13

**NAME:HEMANI SINGH**
**ROLL NO:S119**

# MVLU COLLEGE
# Data analysis with SAS/SSPR/R
# PRACTICAL NO.13





**NAME:HEMANI SINGH**
**ROLL NO:S119**