**Aim:** Applying basic data cleaning functions: handling missing values using na.omit()/replace_na() in R. import dataset.

Code:

```
# R Script: Handling Missing Values (Data Cleaning)

library(dplyr)
library(tidyr) # Contains replace_na()

# 1. IMPORT DATASET

country_df <- read.csv("ESGCountry.csv", na.strings = c("", "NA"))

country_df <- country_df %>%
  rename(alpha_2_code = X2.alpha.code)

print("--- 1. Original Data (Selected Columns, First 6 Rows) ---")
print(head(country_df[, c("Short.Name", "Region", "Income.Group", "Latest.population.census",
"PPP.survey.year")]))

print("--- Count of Missing Values per Column ---")
na_counts <- colSums(is.na(country_df))
print(na_counts[na_counts > 0])

# 2. METHOD A: REMOVE MISSING VALUES (na.omit)

omit_target_cols <- country_df %>%
  select(Short.Name, Region, Income.Group, Lending.category, PPP.survey.year)

clean_omit <- na.omit(omit_target_cols)

print("--- 2. Data after na.omit() ---")
print(paste("Original rows:", nrow(country_df)))
print(paste("Rows remaining:", nrow(clean_omit)))
print(head(clean_omit))


# 3. METHOD B: REPLACE MISSING VALUES (replace_na)

avg_census_year <-
round(mean(as.numeric(as.character(country_df$Latest.population.census)), na.rm = TRUE))
```
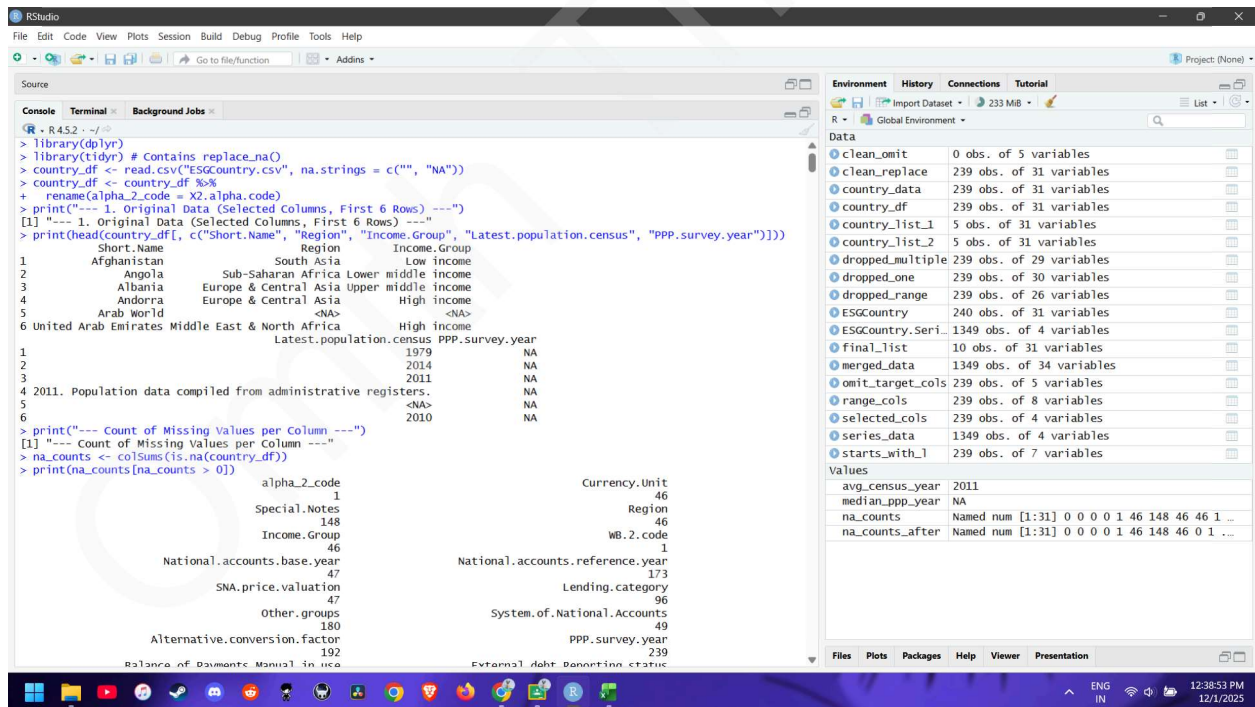
Omith Thilakan
S097

```r
median_ppp_year <- median(country_df$PPP.survey.year, na.rm = TRUE)

clean_replace <- country_df %>%
  replace_na(list(
    Income.Group = "Not Classified",
    Latest.population.census = as.character(avg_census_year),
    PPP.survey.year = median_ppp_year
  ))

print("--- 3. Data after replace_na() ---")
print(clean_replace[is.na(country_df$Income.Group) | is.na(country_df$PPP.survey.year),
             c("Short.Name", "Income.Group", "Latest.population.census", "PPP.survey.year")])
print(head(clean_replace[, c("Short.Name", "Income.Group", "Latest.population.census",
"PPP.survey.year")]))

print("--- Remaining NAs after replacement ---")
na_counts_after <- colSums(is.na(clean_replace))
print(na_counts_after[c("Income.Group", "Latest.population.census", "PPP.survey.year")])
```

Output:



Omith Thilakan
S097

Omith Thilakan
S097

```
225             NA
226             NA
227             NA
228             NA
229             NA
230             NA
231             NA
232             NA
233             NA
234             NA
235             NA
236             NA
237             NA
238             NA
239             NA
> print(head(clean_replace[, c("Short.Name", "Income.Group", "Latest.population.census", "PPP.survey.year")]))
        Short.Name       Income.Group                              Latest.population.census
1       Afghanistan         Low income                                                 1979
2            Angola Lower middle income                                                 2014
3           Albania Upper middle income                                                 2011
4           Andorra        High income 2011. Population data compiled from administrative registers.
5         Arab World     Not Classified                                                 2011
6 United Arab Emirates       High income                                                 2010
  PPP.survey.year
1              NA
2              NA
3              NA
4              NA
5              NA
6              NA
> print("--- Remaining NAs after replacement ---")
[1] "--- Remaining NAs after replacement ---"
> na_counts_after <- colSums(is.na(clean_replace))
> print(na_counts_after[c("Income.Group", "Latest.population.census", "PPP.survey.year")])
            Income.Group Latest.population.census          PPP.survey.year
                       0                        0                      239
>
>
>
>
```

Omith Thilakan
S097