# KDI Foundational DB

Harmonizing Macro-Level Democratic Indicators

Moosa Ghaeb

August 2, 2025

## Project Overview

This document contains the R code for a data pipeline project designed to harmonize several key datasets in comparative political science. The project's objective is to ingest raw data from The World Bank, the V-Dem Institute, and the Quality of Government (QoG) Institute, process and link them deterministically, and load the final, analysis-ready dataset into a PostgreSQL database.

The pipeline is structured across several scripts, each with a distinct purpose, demonstrating a modular and reproducible workflow.

---

## Script 1: Data Ingestion

`R/01_ingest_data.R`

This script handles the "Extract" phase. It downloads all necessary raw data from the three sources and saves them locally.

```
1  # -----------------------------------------------------------------------------
2  # SCRIPT: 01_ingest_and_assess_data.R
3  # PROJECT: KDI Foundational Database
4  #
5  # PURPOSE: 1. Load necessary libraries.
6  #          2. Ingest raw data from three sources.
7  #
8  # -----------------------------------------------------------------------------
9
10
11 if (!require("wbstats")) install.packages("wbstats")
12 if (!require("dplyr")) install.packages("dplyr")
13 if (!require("haven")) install.packages("haven")
14
15 library(wbstats)
16 library(dplyr)
17 library(haven)
18
19
20
21 dir.create("data/vdem", showWarnings = FALSE, recursive = TRUE)
22 dir.create("data/qog", showWarnings = FALSE, recursive = TRUE)
23
```

```
24
25  # World Bank Data (via API)
26  wb_raw_data <- wb_data(
27      indicator = "NY.GDP.PCAP.KD",
28      start_date = 1990,
29      end_date = 2023
30  )
31
32
33  # V-Dem Data
34  vdem_url <- "https://www.v-dem.net/media/datasets/V-Dem-CY-Core-v15_dta.zip"
35  vdem_zip_path <- "data/vdem/V-Dem-CY-Core-v15.zip"
36  vdem_dta_path <- "data/vdem/V-Dem-CY-Core-v15.dta"
37
38  if (!file.exists(vdem_zip_path)) {
39      download.file(vdem_url, destfile = vdem_zip_path, mode = "wb")
40  }
41  if (!file.exists(vdem_dta_path)) {
42      unzip(vdem_zip_path, exdir = "data/vdem")
43  }
44  vdem_raw_data <- read_dta(vdem_dta_path)
45
46
47  # QoG Data
48  qog_url <- "https://www.qogdata.pol.gu.se/data/qog_std_ts_jan25.csv"
49  qog_csv_path <- "data/qog/qog_std_ts_jan25.csv"
50
51  if (!file.exists(qog_csv_path)) {
52      download.file(qog_url, destfile = qog_csv_path, mode = "wb")
53  }
54  qog_raw_data <- read.csv(qog_csv_path)
```

## Script 2: Data Transformation and Merging

`R/02_transform_and_merge.R`

This script performs the "Transform" phase. It cleans each raw dataset, filters for the relevant time period (1990+), standardizes column names, and merges the three sources into a single, coherent data frame using a key-based full join.

```
1   # ----------------------------------------------------------------------------
2   # SCRIPT: 02_transform_and_merge.R
3   # PROJECT: KDI Foundational Database
4   #
5   # PURPOSE: 1. Cleans and standardizes each of the three datasets.
6   #          2. Filters each dataset to the relevant time period (1990+).
7   #          3. Merges the clean datasets into a single master data frame.
8   #
9   # ----------------------------------------------------------------------------
10
11
12  # world bank
13  wb_clean <- wb_raw_data %>%
14      filter(date >= 1990) %>%
15      select(iso3c, country, date, NY.GDP.PCAP.KD) %>%
16      rename(
17          country_iso3 = iso3c,
18          country_name = country,
19          year = date,
20          wb_gdp_pc = NY.GDP.PCAP.KD
21      )
22
```

```r
23
24 # vdem
25 vdem_clean <- vdem_raw_data %>%
26   filter(year >= 1990) %>%
27   select(country_text_id, country_name, year, v2x_polyarchy) %>%
28   rename(
29     country_iso3 = country_text_id,
30     vdem_polyarchy = v2x_polyarchy
31   )
32
33
34 # qog
35 qog_clean <- qog_raw_data %>%
36   filter(year >= 1990) %>%
37   select(ccodealp, cname, year, bci_bci) %>%
38   rename(
39     country_iso3 = ccodealp,
40     country_name = cname,
41     qog_bci = bci_bci
42   )
43
44
45
46 # join
47 merged_data_temp <- vdem_clean %>%
48   full_join(qog_clean, by = c("country_iso3", "year")) %>%
49   full_join(wb_clean, by = c("country_iso3", "year"))
50
51
52 merged_data <- merged_data_temp %>%
53   mutate(country_name = coalesce(country_name.x, country_name.y, country_name))
        %>%
54   select(country_iso3, country_name, year, vdem_polyarchy, qog_bci, wb_gdp_pc)
```

## Script 3: Load Data to PostgreSQL

`R/03_load_to_db.R`

This script handles the "Load" phase. It connects to a PostgreSQL database and writes the final, merged data frame to a table, making it available for analysis.

```r
1  # -----------------------------------------------------------------------------
2  # SCRIPT: 03_load_to_db.R
3  # PROJECT: KDI Foundational Database
4  #
5  # PURPOSE: This script takes the final, merged dataframe from the R
6  #          environment and writes it to the 'harmonized_data' table
7  #          in the PostgreSQL database.
8  #
9  # -----------------------------------------------------------------------------
10
11
12 if (!require("DBI")) install.packages("DBI")
13 if (!require("RPostgres")) install.packages("RPostgres")
14
15 library(DBI)
16 library(RPostgres)
17
18
19 # Establish the connection using credentials stored in .Renviron
20 con <- dbConnect(RPostgres::Postgres(),
21                  dbname = "kdi_db",
```

```
22                    host = "localhost",
23                    port = 5432,
24                    user = "postgres",
25                    password = Sys.getenv("DB_PASSWORD")
26 )
27
28
29 dbWriteTable(con,
30              "harmonized_data",
31              merged_data,
32              overwrite = TRUE,
33              row.names = FALSE)
34
35 dbDisconnect(con)
```

## Script 4: Methodological Extension - Imputation

`R/04_imputation_example.R`

This script demonstrates a methodological extension for handling missing data. It uses linear interpolation to impute missing values in the time-series data for GDP and creates a visualization to illustrate the result.

```
1  # -----------------------------------------------------------------------------
2  # SCRIPT: 04_imputation_example.R
3  #
4  # PURPOSE: linear interpolation.
5  # -----------------------------------------------------------------------------
6
7
8  if (!require("dplyr")) install.packages("dplyr")
9  if (!require("imputeTS")) install.packages("imputeTS")
10 if (!require("ggplot2")) install.packages("ggplot2")
11
12 library(dplyr)
13 library(imputeTS)
14 library(ggplot2)
15
16
17
18 imputed_data <- merged_data %>%
19   group_by(country_iso3) %>%
20   arrange(year) %>%
21   mutate(
22     imputed_gdp = if (sum(!is.na(wb_gdp_pc)) >= 2) {
23       na_interpolation(wb_gdp_pc)
24     } else {
25       wb_gdp_pc
26     }
27   ) %>%
28   ungroup()
29
30
31
32 imputed_data %>%
33   filter(country_iso3 == "POL") %>%
34   ggplot(aes(x = year)) +
35   geom_line(aes(y = imputed_gdp), color = "dodgerblue", linewidth = 1) +
36   geom_point(aes(y = wb_gdp_pc), color = "red", size = 2.5) +
37   labs(
38     title = "GDP per Capita for Poland (Imputed vs. Original)",
```

```
39      subtitle = "Red points are original data. The blue line shows the imputed
        time-series.",
40      x = "Year",
41      y = "GDP per Capita (Constant 2015 US$)",
42      caption = "Source: World Bank data with linear interpolation for missing
        values."
43  ) +
44  theme_minimal()
```