

ALL THE CODES AND RESULTS ARE AVILABLE AT THE REPO:

🌐 [abhi-1408-shek/Innovate_with_GolStats](https://github.com/abhi-1408-shek/Innovate_with_GolStats)

Documentation of Software Tools and Code

Software Tools and Libraries

The analysis was conducted using **Python 3.10**, with the following libraries:

- **pandas (v2.0.3)** for data manipulation, merging, and cleaning.
- **numpy (v1.24.3)** for numerical operations and handling large datasets.
- **openpyxl (v3.1.2)** to parse the Excel layout file (Layout_HCES_2022-23.xls).
- **concurrent.futures** (built-in) for parallel processing of hierarchical data files.
- **plotly (v5.18.0)** and **seaborn (v0.12.2)** for interactive and statistical visualizations.
- **streamlit (v1.28.0)** to deploy the policy-ready dashboard.

Data Processing Workflow

Reading Layout File

The Layout_HCES_2022-23.xls file defined the structure of raw fixed-width text files. Columns were decoded using byte positions and variable names specified in the layout.

```
#CREATING LAYOUT.csv for FURTHER STEPS.
import pandas as pd

layout_data = pd.read_excel("./Documentation/Layout_HCES 2022-23_modified.xlsx", engine="openpyxl")

new_names = ["Slno", "Item", "QSec", "QItem", "QCol", "Length", "Byte_Start_position", "Dash", "Byte_End_position", "Remarks"]
layout_data.columns = new_names

levels = layout_data[layout_data["Item"].astype(str).str.contains("LEVEL", na=False)].index.tolist()

common_id_rows = layout_data.iloc[4:19][["Slno", "Item", "Length"]].dropna(subset=["Length"])

def extract_level_layout(start_row, next_level_start=None):
    if next_level_start is None:
        level_layout = layout_data.iloc[start_row:][["Slno", "Item", "Length"]]
    else:
        level_layout = layout_data.iloc[start_row:next_level_start][["Slno", "Item", "Length"]]
    return level_layout.dropna(subset=["Length"])

combined_layouts = []
```

Parallel Processing of Fixed-Width Files

The 15 hierarchical files were processed in parallel to optimize efficiency.

```

from concurrent.futures import ProcessPoolExecutor

layout = pd.read_csv("./Output/Layout.csv")

print(layout.head())

levels = layout["Level"].unique()

State_list = pd.read_excel("Documentation/tabulation_state_code.xlsx", usecols=["st", "stn"])

State_list["st"] = pd.to_numeric(State_list["st"], errors='coerce')
print(State_list.head())

def read_fwf_level(level):
    file_name = f"./RawData/hces22_lvl_{level:02d}.TXT" #creating csv's

    if not os.path.exists(file_name):
        print(f"File not found: {file_name}, skipping level {level}")
        return None

    current_layout = layout[layout["Level"] == level]

    column_widths = current_layout["Length"].tolist()
    column_names = current_layout["Item"].apply(lambda x: x.replace(" ", "_")).tolist()

    df = pd.read_fwf(file_name, widths=column_widths, names=column_names, dtype=str)

    print(f"Columns in level {level} file:", df.columns)

```

Merging and Cleaning Data

The consolidated dataset was cleaned to handle missing values and categorical responses.

```

for file_path in file_paths:
    try:
        df = pd.read_csv(file_path, dtype=str)
        df.drop_duplicates(inplace=True)

        if 'HH_ID' in df.columns:
            if final_df is None:
                final_df = df
            else:
                final_df = pd.merge(final_df, df, on='HH_ID', how='outer', suffixes=('', '_dup'))

                final_df = final_df.loc[:, ~final_df.columns.duplicated()]
        else:
            print(f"Skipping {file_path} - 'HH_ID' column missing")
    except Exception as e:
        print(f"Error processing {file_path}: {e}")

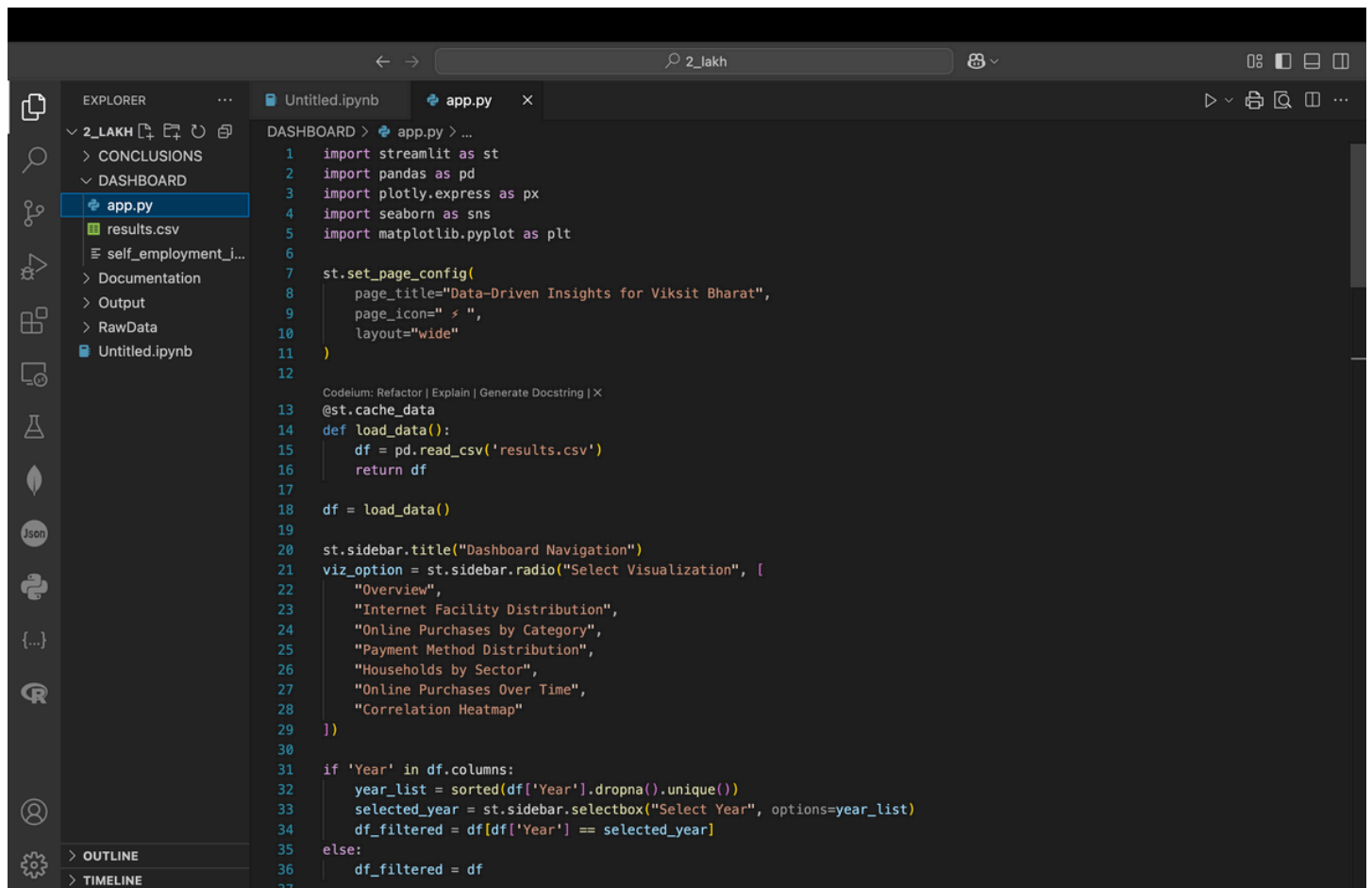
output_file = os.path.join(data_dir, "results.csv")
if final_df is not None:
    final_df.to_csv(output_file, index=False)
    print(f"Consolidation complete. File saved to {output_file}")
else:
    print("No valid data to save.")

```

Data Visualization

Interactive Streamlit Dashboard

The dashboard enables policymakers to explore data dynamically.



Reproducibility

Dependency Installation--->In the app.py file; open console and run:

```
pip install -r requirements.txt
```

Running the Dashboard---> in the same terminal

```
streamlit run app.py
```

Code Repository

All code, datasets, and outputs are available at [Innovate with GoStats](#). The repository includes:

- **Untitled.py**: Script for data preprocessing.
- **app.py**: Streamlit dashboard code.
- **CONCLUSION**: Policy recommendations and reports.

- **DOCUMENTATION:** Layout and state code files.
 - **OUTPUT:** Processed datasets and visualizations.
-

Ethical Compliance

- **Anonymization:** Removed personal identifiers (HH_ID, Person_Srl_No.).
 - **Data Use:** Complied with MoSPI's non-disclosure terms.
-
-

Note: The full code (Untitled.py, app.py) is embedded in the subsequent pages of this document for review.

```
In [ ]: #ABHISHEK SHARMA - Data-Driven Insight

#KINDLY REFER TO THE PROJECT FILE IN THE CONCLUSIONS FOLDER FOR MORE DETAILS
```

```
In [ ]: #CREATING LAYOUT.csv for FURTHER STEPS.
import pandas as pd

layout_data = pd.read_excel("./Documentation/Layout_HCES 2022-23_modified.xlsx")

new_names = ["Slno", "Item", "QSec", "QItem", "QCol", "Length", "Byte_Start"]
layout_data.columns = new_names

levels = layout_data[layout_data["Item"].astype(str).str.contains("LEVEL", rna

common_id_rows = layout_data.iloc[4:19][["Slno", "Item", "Length"]].dropna(s

def extract_level_layout(start_row, next_level_start=None):
    if next_level_start is None:
        level_layout = layout_data.iloc[start_row:][["Slno", "Item", "Length"]]
    else:
        level_layout = layout_data.iloc[start_row:next_level_start][["Slno",
    return level_layout.dropna(subset=["Length"])

combined_layouts = []

for i in range(len(levels)):
    start_row = levels[i] + 3
    next_level_start = levels[i + 1] if i < len(levels) - 1 else len(layout_
    level_layout = extract_level_layout(start_row, next_level_start)

    if i == 0:
        combined_layout = level_layout
    else:
        combined_layout = pd.concat([common_id_rows, level_layout], ignore_i

    combined_layouts.append(combined_layout)

final_combined_layout = pd.concat(combined_layouts, keys=range(1, len(combir

final_combined_layout = final_combined_layout[final_combined_layout["Item"]

final_combined_layout.to_csv("./Output/Layout.csv", index=False)
#SAVING IN OUTPUT FOLDER
```

```
In [ ]: !pip install pandas numpy openpyxl pyreadr multiprocessing
#installing dependencies
```

```

Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (2.0.2)
Requirement already satisfied: openpyxl in /usr/local/lib/python3.11/dist-packages (3.1.5)
Collecting pyreadr
  Downloading pyreadr-0.5.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (1.4 kB)
Collecting multiprocessing
  Downloading multiprocessing-0.70.17-py311-none-any.whl.metadata (7.2 kB)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
Requirement already satisfied: et-xmlfile in /usr/local/lib/python3.11/dist-packages (from openpyxl) (2.0.0)
Collecting dill>=0.3.9 (from multiprocessing)
  Downloading dill-0.3.9-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
Downloading pyreadr-0.5.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (411 kB)
_____ 411.7/411.7 kB 9.1 MB/s eta 0:00
0:00
Downloading multiprocessing-0.70.17-py311-none-any.whl (144 kB)
_____ 144.3/144.3 kB 9.2 MB/s eta 0:00
0:00
Downloading dill-0.3.9-py3-none-any.whl (119 kB)
_____ 119.4/119.4 kB 8.5 MB/s eta 0:00
0:00
Installing collected packages: dill, multiprocessing, pyreadr
Successfully installed dill-0.3.9 multiprocessing-0.70.17 pyreadr-0.5.3

```

```

In [ ]: #DATA PREPROCESSING MULTI-THREADING
import pandas as pd
import numpy as np
import os
from concurrent.futures import ProcessPoolExecutor

layout = pd.read_csv("./Output/Layout.csv")

print(layout.head())

levels = layout["Level"].unique()

State_list = pd.read_excel("Documentation/tabulation_state_code.xlsx", usecols="st")

State_list["st"] = pd.to_numeric(State_list["st"], errors='coerce')
print(State_list.head())

def read_fwf_level(level):
    file_name = f"./RawData/hces22_lvl_{level:02d}.TXT" #creating csv's

```

```

if not os.path.exists(file_name):
    print(f"File not found: {file_name}, skipping level {level}")
    return None

current_layout = layout[layout["Level"] == level]

column_widths = current_layout["Length"].tolist()
column_names = current_layout["Item"].apply(lambda x: x.replace(" ", "_"))

df = pd.read_fwf(file_name, widths=column_widths, names=column_names, dtype=object)

print(f"Columns in level {level} file:", df.columns)

df["Multiplier"] = pd.to_numeric(df.get("Multiplier"), errors="coerce")
df["Weights"] = df["Multiplier"] / 100

fsu_col = next((col for col in df.columns if "FSU" in col), None)
stratum_col = next((col for col in df.columns if "Second" in col), None)
hhld_col = next((col for col in df.columns if "Sample" in col), None)

if fsu_col and stratum_col and hhld_col:
    df["HH_ID"] = df[fsu_col].fillna("").astype(str) + \
        df[stratum_col].fillna("").astype(str) + \
        df[hhld_col].fillna("").astype(str)

if "State" in df.columns:
    df["State"] = pd.to_numeric(df["State"], errors="coerce")
    df = df.merge(State_list, left_on="State", right_on="st", how="left")

return df

num_workers = max(os.cpu_count() - 2, 1)
with ProcessPoolExecutor(max_workers=num_workers) as executor:
    data_frames = list(executor.map(read_fwf_level, levels))

valid_data = [(level, df) for level, df in zip(levels, data_frames) if df is not None]

output_dir = "./Output"
os.makedirs(output_dir, exist_ok=True)

for level, df in valid_data:
    file_prefix = f"level_{level}"
    df.to_csv(os.path.join(output_dir, f"{file_prefix}.csv"), index=False)
    df.to_pickle(os.path.join(output_dir, f"{file_prefix}.pkl"))

del layout, State_list, data_frames

print("Processing complete.")

```

	Level	Slno	Item	Length
0	1	1	Survey Name	4
1	1	2	Year	4
2	1	3	FSU Serial No.	5
3	1	4	Sector	1
4	1	5	State	2

	st	stn
0	28	Andhra Pradesh
1	12	Arunachal Pradesh
2	18	Assam
3	10	Bihar
4	22	Chattisgarh

Columns in level 1 file: Index(['Survey_Name', 'Year', 'FSU_Serial_No.', 'Sector', 'State',

'NSS-Region', 'District', 'Stratum', 'Sub-stratum', 'Panel',
'Sub-sample', 'FOD-Sub-Region', 'Sample_SU_No.',
'Sample_Sub-division_no.', 'Second-stage-stratum_no.',
'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'Survey_Code',
'Reason_for_substitution_Code', 'Multiplier'],

dtype='object')

Columns in level 2 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',

'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n

o.',

'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'Person_Srl_No.',
'Relation_to_head(code)', 'Gender', 'Age(in_years)',
'Marital_Status(code)', 'Highest_educational_level_attained(code)',
'Total_year_of_education_completed',
'Whether_used_internet_from_any_location_during_last_30_days',
'No._of_days_stayed_away_from_home_during_last_30_days',
'No._of_meals_usually_taken_in_a_day',
'No._of_meals_taken_during_last_30_days_from_school_balwadi_etc.',
'No._of_meals_taken_during_last_30_days_from_employer_as\nperquisites

_or_part_of_wage',

'No._of_meals_taken_during_last_30_days_others',
'No._of_meals_taken_during_last_30_days_on_payment',
'No._of_meals_taken_during_last_30_days_at_home',
'Status_of_Member_as_on_revisit',
'FDQ_original_member(generated_field)', 'Multiplier'],

dtype='object')

Columns in level 3 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',

'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n

o.',

'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'HH_Size_(For_FD

Q)',

'Whether_any_household_member_(excluding_those_employed_by_the_househ
old_and_paying_guests)_was_engaged_in_economic_activities_during_last_365_da
ys?_',

'NCO-2015_Code(3-digit)', 'NIC-2008_Code(5-digit)',

'Broad_activities_from_which_maximum_income_was_derived_by_the_househ
old_during_last_365_days',

'Whether_major_source_of_income_from_self-employment_was_from_agricu
lture_or_agricultural_sector',


```

        'Whether_the_major_income_from_regular_wage/salary_earning_from_agricultural_/non-agricultural_sector',
        'Whether_the_major_income_from_casual_labour_was_from_agriculturalsector/_non-agricultural_sector',
        'Household_Type', 'Religion_of_the_head_of_the_household',
        'Social_Group_of_the_head_of_the_household',
        'Does_the_household_own_(owned_&_possessed_or_leased_out)_any_land_(within_the_country)_as_on_the_date_of_survey?',
        'Type_of_land_owned',
        'What_is_the_total_area_of_all_owned_(owned_and_leased_or_leased_out)_land_(within_the_country)_by_the_household_as_on_the_date_of_survey_(are_a_in_acre)?(upto_two_places_of_decimal)',
        'Does_the_household_have_a_dwelling_unit_at_present_place_of_enumeration?__ ',
        'Type_of_Dwelling_Unit',
        'Basic_building_Material_used_for_major_portion_of_the_wall_of_the_dwelling_Unit',
        'Basic_building_Material_used_for_construction_of_the_major_portion_of_the_outer_exposed_part_of_the_roof_of_the_dwelling_unit',
        'Basic_Building_Material_used_for_construction_of_the_major_portion_of_the_floor_of_the_dwelling_Unit`',
        '_Primary_source_of_energy_of_the_household_for_cooking',
        '_Primary_source_of_energy_of_the_household_for_Lighting',
        'Source_of_Drinking_Water_(Last_365_days)',
        'Time_taken_by_the_household_for_a_single_trip_to_reach_the_source_(from_which_most_of_the_drinking_water_is_fetched),_obtain_water_and_back_to_household_(in_Minutes)',
        'Type_of_access_of_the_household_to_latrine',
        'Type_of_latrine_in_which_the_household_has_access',
        'Type_of_ration_card_possessed_by_the_household_as_on_the_date_of_survey',
        'Prevailing_rate_of_rent_in_the_locality_is_available_(FOR_Rural_only)',
        'Benefitted_from_PMGKY_as_on_the_date_on_the_survey',
        'Any_member_of_the_household_of_age_0_-_18_years_died_during_the_period_of_last_5_years_preceding_the_date_of_survey',
        'No_of_members_of_the_household_of_age_0_-_18_years_died_during_the_period_of_last_5_years_preceding_the_date_of_survey?',
        'Multiplier'],
        dtype='object')
Columns in level 4 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State', 'NSS-Region', 'District',
        'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region', 'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_no.',
        'Sample_hhld_No.', 'Questionnaire_No.', 'Level',
        'Whether_the_household_procured_any_item_using_the_ration_card_during_the_last_30_days',
        'Which_item_procured_using_ration_card_during_the_last_30_days_-_Rice',
        'Which_item_procured_using_ration_card_during_the_last_30_days_-_Wheat',
        'Which_item_procured_using_ration_card_during_the_last_30_days_-_Coarse_Grain',
        'Which_item_procured_using_ration_card_during_the_last_30_days_-_Sugar']

```

```

        'Which_item_procured_using_ration_card_during_the_last_30_days-_Pulse
s',
        'Which_item_procured_using_ration_card_during_the_last_30_days-_Edibl
e_Oil',
        'Which_item_procured_using_ration_card_during_the_last_30_days-_Other
_food_item_(including_salt,gram,etc.)',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Groceries',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Milk_&_its_products',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Vegetables',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Fresh_Fruits',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Dry_Fruits',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Egg,fish_&_meat',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Served_processed_food',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Packed_processed_food',
        'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Other_food_items',
        'Ceremony_Pereformed_during_last_30_days',
        'Meals_served_to_non-household_members_during_the_last_30_days',
        'Multiplier'],
dtype='object')
Columns in level 5 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',
        'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
        'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
        'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
        'Consumption_out_of_home_produce-Quantity(0.000)',
        'Consumption_out_of_home_produce-Value(Rs.)',
        'Total_Consumption--Quantity(0.000)', 'Total_Consumption--Value(R
s.)'],
        'Source', 'Multiplier'],
dtype='object')
Columns in level 6 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',
        'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
        'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
        'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
        'Total_Consumption-Quantity(0.000)', 'Total_Consumption-Value(Rs.)',
        'Source', 'Multiplier'],
dtype='object')
Columns in level 7 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',
        'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
        'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
        'Sample_hhld._No.', 'Questionnaire_No.', 'Level',
        'Whether_the_household_procured_kerosene_using_ration_card_during_las

```

```

t_30_days_-_Kerosene',
    'Whether_household_received_subsidy_on_LPG_cylinder_during_the_last_3_months?',
    'If_yes_in_Q4.2.1_Number_of_subsidized_LPG_cylinder_received_during_the_last_3_months_preceding_the_date_of_survey_(number)_',
    'Whether_household_received_free_electricity_during_the_last_30_days',
    'Whether_any_household_member_is_attending/attended_educational_institution_during_last_365_days',
    'If_Code_1_in_Q4.2.3_Number_attending_/attended_Govt._Institution_?',
    'If_Code_1_in_Q4.2.3_Number__attended_/attended_Private_Institution_?',
    '[Checkbox]_If_Code_1_in_Q4.2.3_whether_any_member_of_the_household_received_following_items_free_in_last_365_days:_Textbooks_?',
    'Total_no._of_free_textbooks_received',
    '[Checkbox]_If_Code_1_in_Q4.2.3_whether_any_member_of_the_household_received_following_items_free_in_last_365_days:_Stationary_(pen,notepad_et_c.)_?',
    'Total_no._of_free_stationaries_received',
    '[Checkbox]_If_Code_1_in_Q4.2.3_whether_any_member_of_the_household_received_following_items_free_in_last_365_days:_School_Bag_?',
    'Total_no._of_free_school_bags_received',
    '[Checkbox]_If_Code_1_in_Q4.2.3_whether_any_member_of_the_household_received_following_items_free_in_last_365_days:_Others_?',
    'Total_no._of_free_other_items_received',
    'If_Code_1_in_Q4.2.3_Whether_Any_member_received_reimbursement/waiver_of_school/clg._Fee_during_last_365_days?',
    'If_Code_1_in_Q4.2.6_Number_of_member_received_reimbursement/waiver_?',
    'Is_one_or_more_member_of_the_household_a_benificiary_of_Pradhan_Mantri_Jan_Aarogya_Yojana_(Ayushman_Bharat)_or_any_other_state_specific_public_health_scheme_as_on_the_date_of_survey',
    'If_Code_1_in_Q4.2.7_Number_of_beneficiaries?',
    'Whether_there_was_any_case_of_hospitalization_in_the_household_during_last_365_days?',
    'Whether_one_or_more_member_of_the_household_has_received_benefits_of_medical_treatment_(medical_-_hospitalisation)under_Pradhan_Mantri_Jan_Aarogya_Yojana_Card_(Ayushman_Bharat)_or_any_other_state_specific_public_health_scheme_during_the_last_365_days',
    'If_yes_in_Q4.2.9_number_of_member_received_benifit?',
    'If_yes_in_Q4.2.9_Amount?',
    'Whether_any_online_purchase/payment_has_been_made_during_the_reference_period_to_buy_-_Fuel_&_light',
    'Whether_any_online_purchase/payment_has_been_made_during_the_reference_period_to_buy_-_Toilet_articles_&_other_household_consumables',
    'Whether_any_online_purchase/payment_has_been_made_during_the_reference_period_to_buy_-_Education',
    'Whether_any_online_purchase/payment_has_been_made_during_the_reference_period_to_buy_-_Medicine_&_other_medical_services',
    'Whether_any_online_purchase/payment_has_been_made_during_the_reference_period_to_buy_-_Services_(Travel,_Recharges,_Bill_payment,_Cinema/Theatre,_internet,_etc.)_',
    'Household_has_internet_facility_as_on_the_date_of_the_survey',
    'Multiplier'],
    dtype='object')

```

```

'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
    'Consumption_out_of_home_produce--Quantity(0.000)',
    'Consumption_out_of_home_produce--Value(Rs.)',
    'Total_Consumption--Quantity(0.000)', 'Total_Consumption--Value(R
s.)',
    'Source', 'Multiplier'],
    dtype='object')
Columns in level 9 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
    'Value(Rs.)', 'Multiplier'],
    dtype='object')
Columns in level 10 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld.No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
    'Consumption_out_of_home_produce--Quantity(0.000)',
    'Consumption_out_of_home_produce--Value(Rs.)',
    'Total_Consumption--Quantity(0.000)', 'Total_Consumption--Value(R
s.)',
    'Source', 'Multiplier'],
    dtype='object')
Columns in level 11 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld.No.', 'Questionnaire_No.', 'Level',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Clothing',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Footwear',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Furniture_&
fixtures',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Mobile_hands
et',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Personal_good
s_(laptop/PC,_tablet,_clock,_watch,_spectacles,_contact_lenses,_etc.)',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Goods_for_rec
reation_(TV,_camera,_pen-drive,_musical_instruments,_\netc.)\n',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Cooking_&_oth
er_household_appliances_',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Crockery_&_u
tensils',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Sports_good
s',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Medical equip
ment',
    'Which_online_item_purchased/Paid_during_last_365_days_-_Bedding',

```

```

        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Laptop/PC',
        'Total_number_of_free_Laptop/PC',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Tablet',
        'Total_number_of_free_Tablet',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Mobile',
        'Total_number_of_free_Mobile',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Bicycle',
        'Total_number_of_free_Bicycle',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Motorcycle/Scooty',
        'Total_number_of_free_Motorcycle/Scooty',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Clothing_(Schooling_uniform_etc.)',
        'Total_number_of_free_Clothing_(Schooling_uniform_etc.)',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Footwear_(School_shoe_etc.)',
        'Total_number_of_free_Footwear_(School_shoe_etc.)',
        '[Checkbox]_:Whether_One_or_more_member_of_the_household_received_it
ems_free_of_cost_during_last_365_days:_Other',
        'Total_number_of_free_Other_items',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Television_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Radio_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Laptop/PC_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Mobile_handset_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Bicycle_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Motorcycle,scooter_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Motor_car/jeep/van_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Trucks_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Animal_cart_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Refrigerator_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Washing_machine_',
        'Whether_household_posessed_one_or_more_item_as_on_the_date_of_the_s
urvey-_Air_conditioner/air_cooler_',
        'Type_of_multichannel_television_facility_is_used_by_the_household_as
_on_the_date_of_the_survey',
        'Multiplier'],
        dtype='object')
Columns in level 12 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
        'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
        'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n

```

```

o.',
    'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Item_Code',
    'Quantity_(0.00)', 'Value_(Rs.)', 'Multiplier'],
    dtype='object')
Columns in level 13 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Item_Code_',
    '_First-hand_purchase:number', 'Whether_purchased_on_hire',
    '_First-hand_purchase:Value(Rs.)',
    'Cost_of_repair_&_maintenance_/cost_of_raw_material_and_services_for_
construction_and_repair(Rs.)',
    '2nd-hand_purchase:Number', '_2nd-hand_purchase:Value(Rs.)',
    'Total_expenditure(Rs.)', 'Multiplier'],
    dtype='object')
Columns in level 14 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Section',
    'Item_Code', 'Value_(in_Rs)', 'Multiplier'],
    dtype='object')
Columns in level 15 file: Index(['Year', 'FSU_Serial_No.', 'Sector', 'Stat
e', 'NSS-Region', 'District',
    'Stratum', 'Sub-stratum', 'Panel', 'Sub-sample', 'FOD-Sub-Region',
    'Sample_SU_No.', 'Sample_Sub-division_no.', 'Second-stage-stratum_n
o.',
    'Sample_hhld._No.', 'Questionnaire_No.', 'Level', 'Section',
    'Time_taken_to_canvass_the_questionnaire_(in_minutes)',
    'Household's_usual_consumption_expenditure_in_a_month_(in_Rs.)',
    'Total_expenditure_incurred_on_online_purchase/payment_in_last_30_day
s',
    'Informant_code', 'Response_code', 'Household_size', 'Multiplier'],
    dtype='object')
Processing complete.

```

```

In [ ]: #CONSOLIDATING CSV'S TO A SINGLE FILE i.e Results.csv
import pandas as pd
import os
from glob import glob

data_dir = "./Output"

file_paths = sorted(glob(os.path.join(data_dir, "level_*.csv")))

if not file_paths:
    print("No CSV files found in", data_dir)
else:
    print("Found CSV files:", file_paths)

final_df = None

for file_path in file_paths:

```

```

try:
    df = pd.read_csv(file_path, dtype=str)
    df.drop_duplicates(inplace=True)

    if 'HH_ID' in df.columns:
        if final_df is None:
            final_df = df
        else:
            final_df = pd.merge(final_df, df, on='HH_ID', how='outer', s

            final_df = final_df.loc[:, ~final_df.columns.duplicated()]
    else:
        print(f"Skipping {file_path} - 'HH_ID' column missing")
except Exception as e:
    print(f"Error processing {file_path}: {e}")

output_file = os.path.join(data_dir, "results.csv")
if final_df is not None:
    final_df.to_csv(output_file, index=False)
    print(f"Consolidation complete. File saved to {output_file}")
else:
    print("No valid data to save.")

```

Found CSV files: ['./Output/level_1.csv', './Output/level_10.csv', './Output/level_11.csv', './Output/level_12.csv', './Output/level_13.csv', './Output/level_14.csv', './Output/level_15.csv', './Output/level_2.csv', './Output/level_3.csv', './Output/level_4.csv', './Output/level_5.csv', './Output/level_6.csv', './Output/level_7.csv', './Output/level_8.csv', './Output/level_9.csv']

Consolidation complete. File saved to ./Output/results.csv

```

In [ ]: #checking for duplicates
import pandas as pd

file_path = "./Output/results.csv"
df = pd.read_csv(file_path)

df = df.dropna(axis=1, how="all")

df = df.loc[:, ~df.columns.duplicated()]

print(df.info())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32 entries, 0 to 31
Columns: 180 entries, Survey_Name to Value(Rs.)
dtypes: float64(56), object(124)
memory usage: 45.1+ KB
None

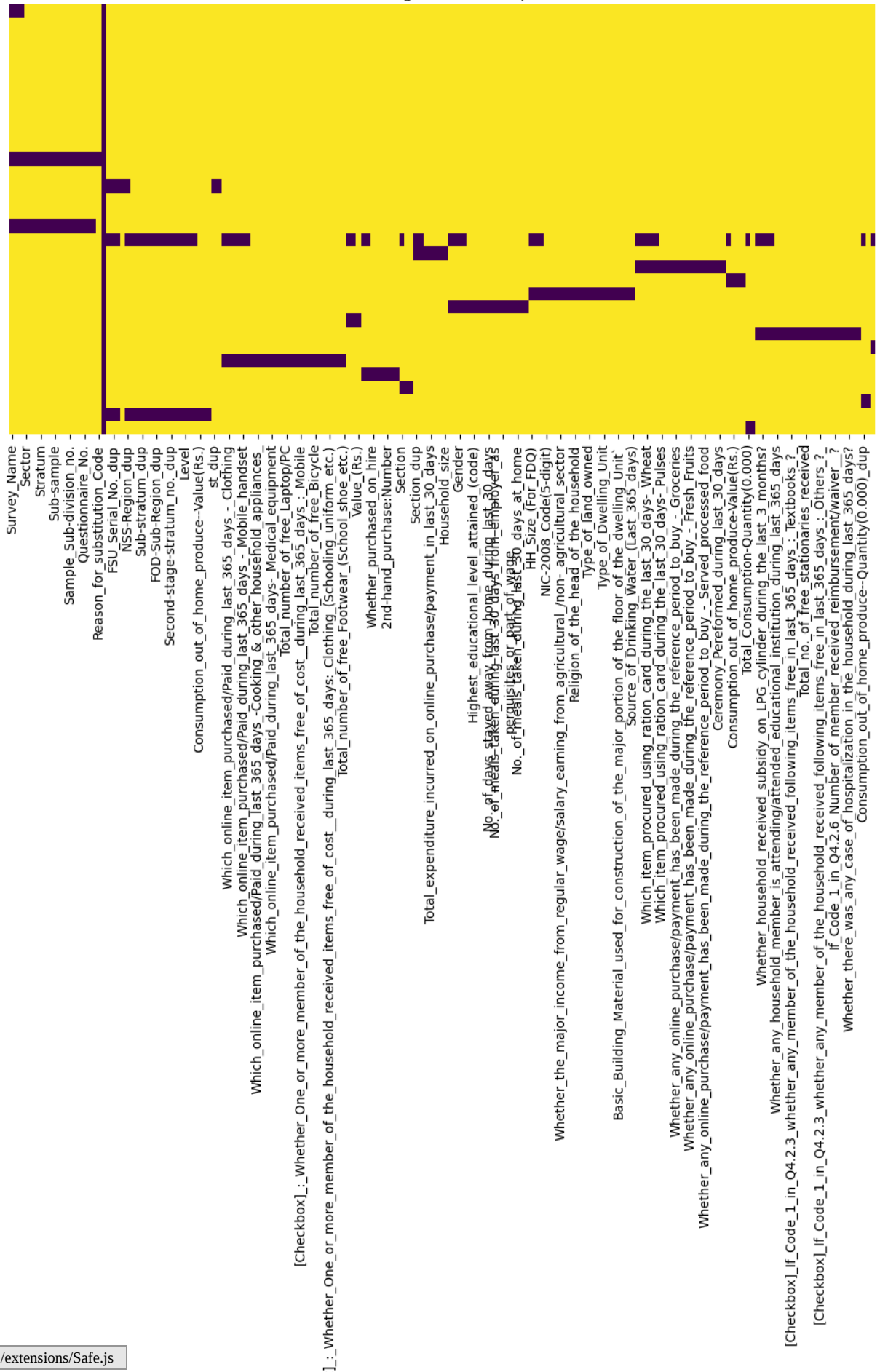
```

```
In [ ]: # Convert columns that should be numeric
for col in df.columns:
    if df[col].dtype == "object":
        try:
            df[col] = pd.to_numeric(df[col])
        except ValueError:
            pass
```

```
In [ ]: #HEATMAP TO CHECK FOR MISSING VALUES BEFORE DATA VISUALIZATION
import matplotlib.pyplot as plt
import seaborn as sns

plt.figure(figsize=(12, 6))
sns.heatmap(df.isnull(), cmap="viridis", cbar=False, yticklabels=False)
plt.title("Missing Data Heatmap")
plt.show()
```


Missing Data Heatmap



```
In [ ]: # CHECKING THE COLUMN VALUES
import pandas as pd
file_path = './Output/results.csv'
df = pd.read_csv(file_path)
df.head()
```

```
Out[ ]:
```

	Survey_Name	Year	FSU_Serial_No.	Sector	State	NSS-Region	District	Stratu
0	size	157	04760	NaN	NaN	NaN	NaN	NaN
1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 219 columns

```
In [16]: # Display column names
print(df.columns)
```

Index(['Survey_Name', 'Year', 'FSU_Serial_No.', 'Sector', 'State',
'NSS-Region', 'District', 'Stratum', 'Sub-stratum', 'Panel',
...
'If_yes_in_Q4.2.9_Amount?',
'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Fuel_&_light',
'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Toilet_articles_&_other_household_consumables',
'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Education',
'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Medicine_&_other_medical_services',
'Whether_any_online_purchase/payment_has_been_made_during_the_referen
ce_period_to_buy_-_Services_(Travel_Recharges_Bill_payment_Cinema/Theatre,
_internet,_etc.)_',
'Household_has_internet_facility_as_on_the_date_of_the_survey',
'Consumption_out_of_home_produce--Quantity(0.000)_dup',
'Consumption_out_of_home_produce--Value(Rs.)_dup', 'Value(Rs.)'],
dtype='object', length=219)

```
In [17]: # Check for columns that might be related to employment
employment_cols = [col for col in df.columns if 'employment' in col.lower()]
print(employment_cols)
```

```
['Whether_major_source_of_income_from_self-_employment_was_from_agricultural  
/_non-_agricultural_sector']
```

```
In [ ]: #SAVED THE "results.csv" file and use it for data visualization further, che
```


DASHBOARD/app.py

```
1 import streamlit as st
2 import pandas as pd
3 import plotly.express as px
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 st.set_page_config(
8     page_title="Data-Driven Insights for Viksit Bharat",
9     page_icon="🇮🇳",
10    layout="wide"
11 )
12
13 @st.cache_data
14 def load_data():
15     df = pd.read_csv('results.csv')
16     return df
17
18 df = load_data()
19
20 st.sidebar.title("Dashboard Navigation")
21 viz_option = st.sidebar.radio("Select Visualization", [
22     "Overview",
23     "Internet Facility Distribution",
24     "Online Purchases by Category",
25     "Payment Method Distribution",
26     "Households by Sector",
27     "Online Purchases Over Time",
28     "Correlation Heatmap"
29 ])
30
31 if 'Year' in df.columns:
32     year_list = sorted(df['Year'].dropna().unique())
33     selected_year = st.sidebar.selectbox("Select Year", options=year_list)
34     df_filtered = df[df['Year'] == selected_year]
35 else:
36     df_filtered = df
37
38 online_purchase_categories = [
39     'Whether_any_online_purchase/payment_has_been_made_during_the_referen-
40 ce_period_to_buy_-_Fuel_&_light',
41     'Whether_any_online_purchase/payment_has_been_made_during_the_referen-
42 ce_period_to_buy_-_Toilet_articles_&_other_household_consumables',
43     'Whether_any_online_purchase/payment_has_been_made_during_the_referen-
44 ce_period_to_buy_-_Education',
45     'Whether_any_online_purchase/payment_has_been_made_during_the_referen-
46 ce_period_to_buy_-_Medicine_&_other_medical_services',
```

```
43     'Whether_any_online_purchase/payment_has_been_made_during_the_referen-
ce_period_to_buy-
_Services_(Travel,_Recharges,_Bill_payment,_Cinema/Theatre,_internet,_etc.)_'
44 ]
45
46 if viz_option == "Overview":
47     st.title("Data-Driven Insights for Viksit Bharat")
48     st.markdown("### Overview")
49     st.dataframe(df.head(10))
50     st.markdown(f"**Total Records:** {df.shape[0]} | **Columns:** {df.shape[1]}")
51
52 elif viz_option == "Internet Facility Distribution":
53     st.title("Internet Facility Distribution in Households")
54     if 'Household_has_internet_facility_as_on_the_date_of_the_survey' in
df_filtered.columns:
55         internet_counts = df_filtered['Household_has_intern-
et_facility_as_on_the_date_of_the_survey'].value_counts()
56         fig = px.pie(
57             names=internet_counts.index,
58             values=internet_counts.values,
59             title="Internet Facility Availability",
60             color_discrete_sequence=px.colors.sequential.RdBu
61         )
62         st.plotly_chart(fig, use_container_width=True)
63
64 elif viz_option == "Online Purchases by Category":
65     st.title("Online Purchases by Category")
66
67     for col in online_purchase_categories:
68         if col in df_filtered.columns:
69             df_filtered[col] = df_filtered[col].map({'Yes': 1, 'No':
0}).fillna(0)
70
71     category_counts = df_filtered[online_purchase_cate-
gories].apply(pd.Series.value_counts).fillna(0)
72
73     if category_counts.shape[1] == 2:
74         category_counts.columns = ['No', 'Yes']
75     elif category_counts.shape[1] == 1:
76         only_response = category_counts.columns[0]
77         if only_response == 0:
78             category_counts = category_counts.rename(columns={0: 'No'})
79             category_counts['Yes'] = 0
80         elif only_response == 1:
81             category_counts = category_counts.rename(columns={1: 'Yes'})
82             category_counts['No'] = 0
83
84     fig = px.bar(
```

```
85     category_counts,
86     barmode='stack',
87     title="Online Purchases Made During the Reference Period",
88     labels={"value": "Number of Households", "index": "Category"},
89     color_discrete_sequence=px.colors.qualitative.Set3
90 )
91 fig.update_layout(xaxis_tickangle=-45)
92 st.plotly_chart(fig, use_container_width=True)
93
94 elif viz_option == "Payment Method Distribution":
95     st.title("Payment Method Distribution")
96     payment_col = 'If_yes_in_Q4.2.9_Amount_?'
97     if payment_col in df_filtered.columns:
98         payment_counts = df_filtered[payment_col].value_counts()
99         fig = px.pie(
100             names=payment_counts.index,
101             values=payment_counts.values,
102             title="Distribution of Online Payment Methods",
103             color_discrete_sequence=px.colors.sequential.Blues
104         )
105         st.plotly_chart(fig, use_container_width=True)
106
107 elif viz_option == "Households by Sector":
108     st.title("Households by Sector")
109     if 'Sector' in df_filtered.columns:
110         sector_counts = df_filtered['Sector'].value_counts().reset_index()
111         sector_counts.columns = ['Sector', 'Count']
112         fig = px.bar(
113             sector_counts,
114             x='Sector',
115             y='Count',
116             title="Households by Sector (Urban vs Rural)",
117             color='Sector',
118             color_discrete_sequence=px.colors.qualitative.Pastel
119         )
120         st.plotly_chart(fig, use_container_width=True)
121
122 elif viz_option == "Online Purchases Over Time":
123     st.title("Online Purchases Over Time")
124     if 'Year' in df.columns:
125         time_trend = df.groupby('Year').size().reset_index(name='Purchases')
126         fig = px.line(
127             time_trend,
128             x='Year',
129             y='Purchases',
130             markers=True,
131             title="Trend of Online Purchases Over the Years",
132             color_discrete_sequence=['#17becf']
```

```
133         )
134         st.plotly_chart(fig, use_container_width=True)
135
136     elif viz_option == "Correlation Heatmap":
137         st.title("Correlation Between Purchase Categories")
138         df_corr = df_filtered.copy()
139         for col in online_purchase_categories:
140             if col in df_corr.columns:
141                 df_corr[col] = pd.to_numeric(df_corr[col].map({'Yes': 1, 'No': 0}),
errors='coerce')
142         corr_matrix = df_corr[online_purchase_categories].corr()
143         fig, ax = plt.subplots(figsize=(10, 8))
144         sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f', ax=ax)
145         st.pyplot(fig)
146
147     st.sidebar.markdown("----")
148     st.sidebar.markdown("Download Cleaned Data")
149     st.sidebar.download_button(
150         label="Download CSV",
151         data=df.to_csv(index=False).encode('utf-8'),
152         file_name="online_purchase_data.csv",
153         mime="text/csv"
154     )
155
156     st.sidebar.text("Created by Abhishek Sharma")
157
```