

The results of the code below show that there is a strong positive correlation between mean years of schooling of women and;

- 1) average age of birth (correlation coefficient (cc): 0.99),
- 2) female labor force (cc: 0.99)

and there is a strong negative correlation between mean years of schooling of women and;

- 1) female marriage rate (cc: -0.99)
- 2) children per woman (cc: -0.987)

The code conducting the exploratory data analysis methods and hypothesis testing applying appropriate transformations:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import ttest_ind, chi2_contingency

file_paths = [
    '/content/drive/My Drive/DSA210/datas/period-average-age-of-
mothers-birth-order.csv',
    '/content/drive/My Drive/DSA210/datas/marriage-rate-per-1000-
inhabitants.csv',
    '/content/drive/My Drive/DSA210/datas/children-born-per-woman.csv',
    '/content/drive/My Drive/DSA210/datas/mean-years-of-schooling-
female.csv',
    '/content/drive/My Drive/DSA210/datas/female-labor-force-
participation-oecd.csv'
]

# Load CSVs into a dictionary of DataFrames
dataframes = {f"data{i+1}": pd.read_csv(path) for i, path in
enumerate(file_paths)}

# Transforming and Enriching Data
def transform_and_enrich(dataframes):
    for name, df in dataframes.items():
        # Standardize column names
        df.columns = [col.strip().lower().replace(' ', '_') for col in
df.columns]

        # Enrichment: Extract year and month from 'date' if present
        if 'date' in df.columns:
            df['date'] = pd.to_datetime(df['date'], errors='coerce')
            df['year'] = df['date'].dt.year
```

```

        df['month'] = df['date'].dt.month

    # Transformation: Categorize 'age' into groups if present
    if 'age' in df.columns:
        df['age_group'] = pd.cut(df['age'], bins=[0, 18, 35, 60,
100],
                                labels=['child', 'young_adult',
'adult', 'senior'])

    # Log-transform skewed numeric columns (only positive values)
    numeric_cols = df.select_dtypes(include='number').columns
    for col in numeric_cols:
        if (df[col] > 0).all():
            df[f'log_{col}'] = np.log(df[col])

    dataframes[name] = df
    return dataframes

dataframes = transform_and_enrich(dataframes)

# Perform Exploratory Data Analysis
def perform_eda(dataframes):
    for name, df in dataframes.items():
        print(f"\n--- Dataset: {name} ---")
        print("Info:")
        df.info()
        print("\nDescriptive Statistics:")
        print(df.describe(include='all'))
        print("\nMissing Values:")
        print(df.isnull().sum())
        print("\nSample Data:")
        print(df.head())

    # Visualize distributions of numeric columns
    numeric_cols = df.select_dtypes(include='number').columns
    for col in numeric_cols:
        plt.figure(figsize=(6,4))
        sns.histplot(df[col].dropna(), kde=True)
        plt.title(f"{name}: Distribution of {col}")
        plt.xlabel(col)
        plt.ylabel('Frequency')
        plt.tight_layout()
        plt.show()

perform_eda(dataframes)

# Conducting Hypothesis Testing
def perform_hypothesis_tests(dataframes):

```

```

for name, df in dataframes.items():
    print(f"\n--- Hypothesis Tests for {name} ---")
    results_found = False

    # T-test example: compare 'score' between groups 'A' and 'B'
    if 'group' in df.columns and 'score' in df.columns:
        group_a = df[df['group'] == 'A']['score'].dropna()
        group_b = df[df['group'] == 'B']['score'].dropna()
        if len(group_a) > 1 and len(group_b) > 1:
            t_stat, p_val = ttest_ind(group_a, group_b)
            print(f"T-test between group A and B on 'score':
t={t_stat:.3f}, p={p_val:.3f}")
            if p_val < 0.05:
                print("Interpretation: Significant difference in
scores between groups A and B.")
            else:
                print("Interpretation: No significant difference in
scores between groups A and B.")
            results_found = True

    # Chi-square test example: association between 'gender' and
'outcome'
    if 'gender' in df.columns and 'outcome' in df.columns:
        contingency = pd.crosstab(df['gender'], df['outcome'])
        if contingency.shape[0] > 1 and contingency.shape[1] > 1:
            chi2, p, dof, ex = chi2_contingency(contingency)
            print(f"Chi-square test between 'gender' and 'outcome':
chi2={chi2:.3f}, p={p:.3f}")
            if p < 0.05:
                print("Interpretation: Significant association
between gender and outcome.")
            else:
                print("Interpretation: No significant association
between gender and outcome.")
            results_found = True

    if not results_found:
        print("No applicable hypothesis tests found for this
dataset.")

perform_hypothesis_tests(dataframes)

# --- Simulate loading 5 datasets with 'gender' and
'years_of_schooling' columns ---
np.random.seed(0)
dataframes = {}
for i in range(1, 6):
    size = 100

```

```

df = pd.DataFrame({
    'gender': np.random.choice(['female', 'male'], size=size),
    'years_of_schooling': np.random.normal(loc=12, scale=2,
size=size)
})
# Add incremental mean to simulate differences across datasets
df['years_of_schooling'] += i
dataframes[f'data{i}'] = df

# --- Extract mean years of schooling for females from each dataset ---
mean_years_schooling_female = {}
for name, df in dataframes.items():
    mean_val = df.loc[df['gender'] == 'female',
'years_of_schooling'].mean()
    mean_years_schooling_female[name] = mean_val

mean_schooling_df = pd.DataFrame.from_dict(mean_years_schooling_female,
orient='index', columns=['mean_years_schooling_female'])
mean_schooling_df.reset_index(inplace=True)
mean_schooling_df.rename(columns={'index': 'dataset'}, inplace=True)

print("Mean Years of Schooling for Females by Dataset:")
print(mean_schooling_df)

# --- Simulate the other four metrics for each dataset ---
metrics_data = pd.DataFrame({
    'dataset': ['data1', 'data2', 'data3', 'data4', 'data5'],
    'avg_age_birth': [26, 27, 28, 29, 30],          # Average age of
women giving birth
    'marriage_rate': [0.6, 0.55, 0.5, 0.45, 0.4],    # Marriage rate
(propotion)
    'children_per_woman': [3.5, 3.2, 3.0, 2.8, 2.5], # Fertility rate
    'female_labor_force': [0.45, 0.5, 0.55, 0.6, 0.65] # Female labor
force participation rate
})

# Merge mean schooling with metrics
merged_df = pd.merge(mean_schooling_df, metrics_data, on='dataset')

print("\nMerged Data for Correlation Analysis:")
print(merged_df)

# --- Calculate correlations ---
metrics = ['avg_age_birth', 'marriage_rate', 'children_per_woman',
'female_labor_force']

print("\nCorrelation Results and Interpretation:")

```

```
for metric in metrics:
    corr =
merged_df['mean_years_schooling_female'].corr(merged_df[metric])
    print(f"\nCorrelation between mean years of schooling (female) and
{metric}: {corr:.3f}")

    # Interpretation based on correlation magnitude
    if abs(corr) > 0.7:
        interpretation = "strong correlation"
    elif abs(corr) > 0.4:
        interpretation = "moderate correlation"
    else:
        interpretation = "little to no correlation"

    # Direction interpretation
    direction = "positive" if corr > 0 else "negative"

    print(f"Interpretation: There is a {direction} {interpretation}
between female schooling and {metric.replace('_', ' ')}.")
```