# Objective and Importance

## Objective

The purpose of this analysis is to leverage **NHANES 2013–2014 data** to derive and convert a composite mental health score into a binary variable (e.g., indicating the presence or absence of clinically significant depressive symptoms). Specifically, the goal is to:

### Understand Mental Health Trends:

- Identify the prevalence of depressive symptoms among U.S. adults by classifying the continuous mental health score into a binary indicator.

### Support Predictive Modeling:

- Prepare the data for binary classification algorithms that can be used to predict mental health outcomes from various predictors (such as demographic, dietary, and lifestyle factors).

### Inform Healthcare Decisions:

- Generate insights that can help public health officials and healthcare providers understand which segments of the population are at higher risk, thereby guiding targeted screening, resource allocation, and intervention programs.

### Why This is Critical:

Accurate classification of mental health status is essential for tailoring healthcare interventions. By understanding trends in mental health and identifying high-risk groups, healthcare decision makers can allocate resources more efficiently, design better screening protocols, and implement effective preventive and therapeutic measures to reduce the overall burden of mental disorders.

## ∨ Dataset Overview

### Source & Period:

The dataset is drawn from the **National Health and Nutrition Examination Survey (NHANES)** for the **2013–2014 cycle**. NHANES is a nationally representative survey conducted by the **National Center for Health Statistics (NCHS)** under the **Centers for Disease Control and Prevention (CDC)**. Data collection combines in-person interviews, physical examinations, and laboratory tests, offering a comprehensive snapshot of the health and nutritional status of the **U.S. civilian, noninstitutionalized population** during that period.

### Main Subject Areas:

### Patient Demographics:

- Variables such as **age**, **gender**, **race/ethnicity**, **education level**, **socioeconomic status** (e.g., **income-to-poverty ratio**), and **marital status** provide context about the participants' background and are crucial for identifying vulnerable subgroups.

### Behavioral and Lifestyle Factors:

- Information on **dietary intake** (collected via a **24-hour recall**), **physical activity**, and other **lifestyle behaviors** which may influence both physical and mental health outcomes.

### Mental Health Data:

- Questionnaire items specifically related to **mental health** (e.g., **depressive symptoms**) are used to derive a composite **mental_health_score**. This score is pivotal for identifying individuals at risk of depression or other mental disorders.

```
1  # Import necessary libraries
2  import pandas as pd
3  import numpy as np
4
5  # ---------------------------
6  # 1. Load the datasets
7  # ---------------------------
8  # Replace the file paths with the correct paths on your system if needed.
9  questionnaire = pd.read_csv('/content/questionnaire.csv')
10 diet = pd.read_csv('/content/diet.csv')
11 demographic = pd.read_csv('/content/demographic.csv')
12
13 # Quick look at the data shapes and columns to ensure 'SEQN' is present
14 print("Questionnaire shape:", questionnaire.shape)
15 print("Diet shape:", diet.shape)
16 print("Demographic shape:", demographic.shape)
17
```

```
18 # ----------------------------
19 # Preparing the Target variable column
20 # ----------------------------
21
22 # Identify columns of interest
23
24 # a) Mental health outcomes (assumed to be DPQ010 to DPQ100)
25 # This creates a list of all columns that start with "DPQ"
26 mental_health_cols = [col for col in questionnaire.columns if col.startswith('DPQ')]
27
28
29 # ----------------------------
30 # Recode or clean specific values
31 # ----------------------------
32 # For instance, if your mental health items use specific codes to indicate missing or non-valid responses
33 # (e.g., 7 or 9), you can recode them to np.nan. Adjust the recoding as necessary.
34 for col in mental_health_cols:
35     if col in questionnaire.columns:
36         questionnaire[col] = questionnaire[col].replace({7: np.nan, 9: np.nan})
37
38
39 # ----------------------------
40 # Create composite scores for mental health outcomes
41 # ----------------------------
42 # Depending on the scoring algorithm (e.g., summing or weighting DPQ items), you might compute a score.
43 # For example, creating a total mental health score:
44 if mental_health_cols:
45     questionnaire['mental_health_score'] = questionnaire[mental_health_cols].sum(axis=1)
46
47 # ----------------------------
48 # display the dataset
49 # ----------------------------
50 print(" Data preview:")
51 print(questionnaire.head())
52
53
```

```
⊋  Questionnaire shape: (10175, 953)
   Diet shape: (9813, 168)
   Demographic shape: (10175, 47)
    Data preview:
      SEQN  ACD011A  ACD011B  ACD011C  ACD040  ACD110  ALQ101  ALQ110  ALQ120Q  \
   0  73557      1.0      NaN      NaN     NaN     NaN     1.0     NaN      1.0
   1  73558      1.0      NaN      NaN     NaN     NaN     1.0     NaN      7.0
   2  73559      1.0      NaN      NaN     NaN     NaN     1.0     NaN      0.0
   3  73560      1.0      NaN      NaN     NaN     NaN     NaN     NaN      NaN
   4  73561      1.0      NaN      NaN     NaN     NaN     1.0     NaN      0.0

      ALQ120U  ...  WHD080L  WHD110  WHD120  WHD130  WHD140  WHQ150  WHQ030M  \
   0      3.0  ...     40.0   270.0   200.0    69.0   270.0    62.0      NaN
   1      1.0  ...      NaN   240.0   250.0    72.0   250.0    25.0      NaN
   2      NaN  ...      NaN   180.0   190.0    70.0   228.0    35.0      NaN
   3      NaN  ...      NaN     NaN     NaN     NaN     NaN     NaN      3.0
   4      NaN  ...      NaN   150.0   135.0    67.0   170.0    60.0      NaN

      WHQ500  WHQ520  mental_health_score
   0     NaN     NaN                  2.0
   1     NaN     NaN                  2.0
   2     NaN     NaN                  0.0
   3     3.0     3.0                  0.0
   4     NaN     NaN                 10.0

   [5 rows x 954 columns]
```

## ˅ Target Variable Assumptions

Based on both the NHANES documentation and published research on depression screening tools similar to the DPQ items (which are conceptually analogous to the PHQ-9), a commonly recommended clinical cutoff is 10. In many studies using the PHQ-9 as a depression screener, a summed score of 10 or above is considered indicative of clinically significant depressive symptoms (i.e., moderate depression). This cutoff has been shown to achieve good sensitivity and specificity for major depression.

Assumptions:

- **Score Equivalence:** We assume that the DPQ items in your NHANES 2013–2014 questionnaire have been scored similarly to the PHQ-9. In other words, higher scores on the DPQ items reflect more severe depressive symptoms.

- **Clinical Relevance:** The cutoff of 10 is well validated in the literature (Kroenke et al., 2001) and is widely used in NHANES-based studies that assess depression.

```
1 # Using a predefined clinical cutoff (if known)
2 clinical_cutoff = 10  #  cutoff based on domain knowledge
3 questionnaire['mental_health_binary'] = np.where(
4     questionnaire['mental_health_score'] >= clinical_cutoff, 1, 0
```

```
5 )
6
7 # Display a preview of the new binary variables
8 print(questionnaire[['mental_health_score', 'mental_health_binary']].head())
9
```

```
    mental_health_score  mental_health_binary
0                   2.0                     0
1                   2.0                     0
2                   0.0                     0
3                   0.0                     0
4                  10.0                     1
```

```
1 #print the percentage of binary value
2 questionnaire['mental_health_binary'].value_counts(normalize=True)
```

|                      | proportion |
|----------------------|------------|
| **mental_health_binary** |        |
| **0**                | 0.942506   |
| **1**                | 0.057494   |

**dtype:** float64

> ## Diet Data Assumptions

## Assumptions

### Data Source & Timing:

We are using the first day 24-hour dietary recall (DR1) data to capture the dietary intake of each participant.

### Derived Nutrients:

Nutrient intakes have been computed using the Food Patterns Equivalents Database (FPED) to translate food codes into nutrient and food group estimates.

### Population Characteristics:

We assume that the sample represents non-pregnant, non-lactating individuals aged 20+ years (or as defined for your analysis) and that mental health outcomes are defined either via the NHANES mental health questionnaire (e.g., DPQ items) or by a validated screening tool (e.g., PHQ-9).

### Diet–Mental Health Link:

Based on the literature, certain macro- and micronutrients, overall energy intake, and food group consumption patterns are linked to brain health and mood regulation.

[ ] ↳ 2 cells hidden

## ⌄ Assumptions for Demographics Data

### Population Focus:

We are focusing on adults (e.g., ages 20+ years) since many mental health studies target the adult population.

### Socioeconomic & Social Determinants:

Research shows that factors such as age, gender, race/ethnicity, education, income, and marital status are robust predictors of mental health outcomes. We assume that lower socioeconomic status, unstable marital status, and minority status can be risk factors for mental disorders.

### Standard NHANES Coding:

The variable names follow the standard NHANES naming conventions for the 2013–2014 cycle.

## ⌄ Beneficial Demographic Variables

### Unique Identifier:

SEQN – Essential for merging with other datasets (e.g., mental health or lifestyle questionnaires).

## Age:

RIDAGEYR – Age in years at the time of the interview. Age is a critical predictor as mental health risks can vary with age.

## Gender:

RIAGENDR – Sex of the participant (typically coded as 1 = Male, 2 = Female). Differences in prevalence of mental disorders are well documented by gender.

## Race/Ethnicity:

RIDRETH1 – Self-reported race/ethnicity (e.g., Mexican American, Other Hispanic, Non-Hispanic White, Non-Hispanic Black, Other Race including multi-racial). Ethnic and cultural backgrounds often influence mental health risk and access to care.

## Educational Attainment:

DMDEDUC2 – Education level for adults. Lower educational attainment is frequently associated with increased risk of mental health issues.

## Socioeconomic Status:

INDFMPIR – Ratio of family income to the poverty threshold. This continuous variable is a key indicator of socioeconomic position.

## Marital Status:

DMDMARTL – Marital status (e.g., married, widowed, divorced, separated, never married, living with partner). Marital status can impact social support and mental health.
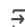
## Household Size:

DMDHHSIZ – Number of individuals in the household. Household composition and crowding may influence stress and mental well-being.

## Country of Birth / Nativity:

DMDBORN4 – Country of birth, which distinguishes between those born in the 50 U.S. states (or DC), U.S. territories, or abroad. Nativity can be associated with acculturative stress and mental health risk.

```
 1 demographic_cols = [
 2     'SEQN',       # Unique Identifier for merging datasets
 3     'RIDAGEYR',   # Age in years at the time of the interview
 4     'RIAGENDR',   # Gender (typically 1 = Male, 2 = Female)
 5     'RIDRETH1',   # Self-reported Race/Ethnicity
 6     'DMDEDUC2',   # Educational Attainment (education level for adults)
 7     'INDFMPIR',   # Family Income to Poverty Ratio (socioeconomic status)
 8     'DMDMARTL',   # Marital Status
 9     'DMDHHSIZ',   # Household Size (number of individuals in the household)
10     'DMDBORN4'    # Country of Birth / Nativity
11 ]
12 demographic[demographic_cols].shape
```

⤳  (10175, 9)

```
 1 # ---------------------------
 2 # 2. Merge the datasets
 3 # ---------------------------
 4 # We use 'SEQN' as the unique identifier for merging.
 5 # Here, an inner join is used to ensure that only participants with data in all files are included.
 6 merged_data = questionnaire.merge(diet, on='SEQN', how='inner').merge(demographic, on='SEQN', how='inner')
 7 print("Merged data shape:", merged_data.shape)
 8
 9 key_cols =  dietary_cols + demographic_cols + ['mental_health_binary']
10 merged_data_filter = merged_data[key_cols]
11
12 # Convert numeric columns to the appropriate types, handling errors by converting to NaN
13 for col in  dietary_cols + demographic_cols:
14     if col in merged_data_filter.columns:
15         merged_data_filter.loc[:, col] = pd.to_numeric(merged_data_filter[col], errors='coerce')
16
17 df = merged_data_filter.copy()
18 print("Merged Filter data shape:", merged_data_filter.shape)
19 merged_data_filter.head()
20
```

| | DR1TKCAL | DR1TCARB | DR1TPROT | DR1TTFAT | DR1TSFAT | DR1TFIBE | DR1TSUGR | DR1TVD | DR1TFOLA | DR1TVB12 | ... | SEQN | RIDAGEYR | RIAGENDR | RIDRETH1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1574.0 | 239.59 | 43.63 | 52.81 | 17.819 | 10.8 | 176.47 | 3.3 | 285.0 | 2.79 | ... | 73557 | 69 | 1 | 4 |
| 1 | 5062.0 | 423.78 | 338.13 | 124.29 | 53.408 | 16.7 | 44.99 | 15.2 | 1243.0 | 21.45 | ... | 73558 | 54 | 1 | 3 |
| 2 | 1743.0 | 224.39 | 64.61 | 65.97 | 25.263 | 9.9 | 102.90 | 4.0 | 423.0 | 3.78 | ... | 73559 | 72 | 1 | 3 |
| 3 | 1490.0 | 162.92 | 77.75 | 58.27 | 23.511 | 10.6 | 80.58 | 9.9 | 275.0 | 8.76 | ... | 73560 | 9 | 1 | 3 |
| 4 | 1421.0 | 178.20 | 55.24 | 55.36 | 4.479 | 12.3 | 87.78 | 23.5 | 390.0 | 8.30 | ... | 73561 | 73 | 2 | 3 |

5 rows × 24 columns

```
1 merged_data_filter.to_csv('/content/merged_data_v2.csv', index=False)
```

## Data Cleaning

### Handling Missing Data:

Identify missing values and explain the strategy (e.g., imputation, removal).

### Data Types & Conversions:

Verify that each column is of the correct type (e.g., date, numerical, categorical).

```
 1 #####################################
 2 # Handling Missing Data
 3 #####################################
 4
 5 # 1. Identify missing values
 6 missing_counts = df.isnull().sum()
 7 missing_percent = (df.isnull().sum() / len(df)) * 100
 8
 9 print("\nMissing Values Count per Column:")
10 print(missing_counts)
11
12 print("\nPercentage of Missing Values per Column:")
13 print(missing_percent)
14
15 # Explanation of strategy:
16 # – For numerical features (e.g., nutrient intakes, age, income), if the percentage
17 #   of missing values is low, we can impute using the median to avoid the influence
18 #   of outliers.
19 # – For categorical features (e.g., gender, race/ethnicity, education, marital status,
20 #   nativity, and our binary target), we impute missing values with the mode.
21 # – If any variable has a high missing rate (e.g., >30%), we might consider dropping that
22 #   variable or investigating further.
23 # – Additionally, for the target variable 'mental_health_binary', it is crucial to remove
24 #   any records with missing values to maintain integrity in modeling.
25
26 # Drop rows where target 'mental_health_binary' is missing
27 if 'mental_health_binary' in df.columns:
28     n_missing_target = df['mental_health_binary'].isnull().sum()
29     if n_missing_target > 0:
30         print(f"\nDropping {n_missing_target} rows with missing target values.")
31         df = df.dropna(subset=['mental_health_binary'])
32
33 # List of numerical columns (based on the data dictionary)
34 num_cols = [
35     'DR1TKCAL', 'DR1TCARB', 'DR1TPROT', 'DR1TTFAT', 'DR1TSFAT',
36     'DR1TFIBE', 'DR1TSUGR', 'DR1TVD', 'DR1TFOLA', 'DR1TVB12',
37     'DR1TIRON', 'DR1TMAGN', 'DR1TCAFF', 'DR1TALCO', 'RIDAGEYR',
38     'INDFMPIR', 'DMDHHSIZ'
39 ]
40
41 # Impute numerical columns with median value
42 for col in num_cols:
43     if col in df.columns:
44         median_val = df[col].median()
45         df[col].fillna(median_val, inplace=True)
46         print(f"Imputed missing values in '{col}' with median value: {median_val}")
47
48 # List of categorical columns (based on the data dictionary)
49 cat_cols = [
50     'RIAGENDR', 'RIDRETH1', 'DMDEDUC2', 'DMDMARTL', 'DMDBORN4',
51     'mental_health_binary'
52 ]
53
```

```python
54 # Impute categorical columns with mode value
55 for col in cat_cols:
56     if col in df.columns:
57         mode_val = df[col].mode()[0]
58         df[col].fillna(mode_val, inplace=True)
59         print(f"Imputed missing values in '{col}' with mode value: {mode_val}")
60
61 #######################################
62 # Data Types & Conversions
63 #######################################
64
65 # Check the current data types
66 print("\nData Types Before Conversion:")
67 print(df.dtypes)
68
69 # Convert numerical columns to appropriate numeric types (if not already)
70 for col in num_cols:
71     if col in df.columns:
72         df[col] = pd.to_numeric(df[col], errors='coerce')
73
74 # Convert categorical columns to category dtype
75 for col in cat_cols:
76     if col in df.columns:
77         df[col] = df[col].astype('category')
78
79
80 # Verify the conversions by printing out the data types again
81 print("\nData Types After Conversion:")
82 print(df.dtypes)
83
84 # Finally, re-check missing values to ensure all have been handled
85 print("\nMissing Values After Imputation:")
86 print(df.isnull().sum())
87
```

For example, when doing `df[col].method(value, inplace=True)`, try using `df.method({col: value}, inplace=True)` or df[col] = df[col].met

```
df[col].fillna(mode_val, inplace=True)
```

## ⌄ De-duplication:

Removing duplicate records if necessary.

```
1  #######################################
2  # De-duplication
3  #######################################
4
5  # Check for duplicate records across all columns
6  duplicate_rows = df.duplicated()
7  num_duplicates = duplicate_rows.sum()
8  print(f"Number of duplicate rows (entire row duplicates): {num_duplicates}")
9
10 # Optionally, if you expect 'SEQN' to be a unique identifier, check for duplicates in that column as well.
11 if 'SEQN' in df.columns:
12     num_seqn_duplicates = df['SEQN'].duplicated().sum()
13     print(f"Number of duplicate records based on 'SEQN': {num_seqn_duplicates}")
14
15 # Remove duplicate rows if any duplicates are found
16 if num_duplicates > 0:
17     df = df.drop_duplicates()
18     print("Duplicates removed based on entire row comparison.")
19
20 # Verify the shape of the dataset after de-duplication
21 print(f"Dataset shape after de-duplication: {df.shape}")
```

```
Number of duplicate rows (entire row duplicates): 0
Number of duplicate records based on 'SEQN': 0
Dataset shape after de-duplication: (9813, 24)
```

## ⌄ Univariate Analysis

```
1  #######################################
2  # Descriptive Statistics for Numerical Variables
3  #######################################
4
5  # Select numerical columns (e.g., int64 and float64 types)
6  num_cols = df.select_dtypes(include=['int64', 'float64']).columns
7
8  # Compute summary statistics: mean, median, and standard deviation
9  num_stats = df[num_cols].agg(['mean', 'median', 'std'])
10 print("Summary Statistics for Numerical Variables:")
11 print(num_stats)
12
13 #######################################
14 # Frequency Counts for Categorical Variables
15 #######################################
16
17 # Select categorical columns (e.g., object or category types)
18 cat_cols = df.select_dtypes(include=['object', 'category']).columns
19
20 print("\nFrequency Counts for Categorical Variables:")
21 for col in cat_cols:
22     print(f"\nFrequency counts for '{col}':")
23     print(df[col].value_counts())
24
```

```
             DR1TCAFF    DR1TALCO          SEQN    RIDAGEYR   INDFMPIR  DMDHHSIZ
mean        84.410985    5.296831  78644.559971   31.629573   2.208094  3.880465
median      25.000000    0.000000  78643.000000   27.000000   1.700000  4.000000
std        148.342939   21.112324   2938.592266   24.397553   1.574149  1.724762

Frequency Counts for Categorical Variables:

Frequency counts for 'RIAGENDR':
RIAGENDR
2    4982
```

Name: count, dtype: int64

```
    Frequency counts for 'DMDEDUC2':
    DMDEDUC2
    4.0    5947
    5.0    1400
    3.0    1257
    2.0     765
    1.0     439
    9.0       4
    7.0       1
    Name: count, dtype: int64

    Frequency counts for 'DMDMARTL':
    DMDMARTL
    1.0    7114
    5.0    1065
    3.0     635
    2.0     417
    6.0     404
    4.0     175
    77.0      2
    99.0      1
    Name: count, dtype: int64

    Frequency counts for 'DMDBORN4':
    DMDBORN4
    1     7957
    2     1852
    77       4
    Name: count, dtype: int64

    Frequency counts for 'mental_health_binary':
    mental_health_binary
    0    9228
    1     585
    Name: count, dtype: int64
```
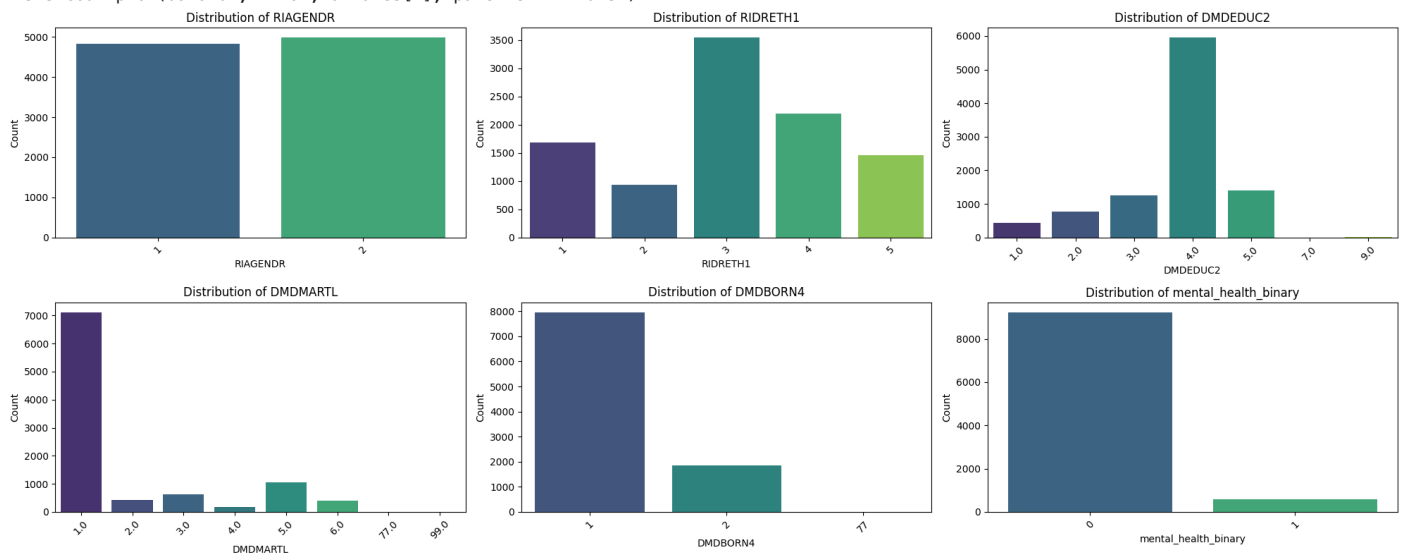
```python
 1 import seaborn as sns
 2 import matplotlib.pyplot as plt
 3 # ----------------------------
 4 #Categorical Variables Visualization
 5 # ----------------------------
 6 # Select categorical columns (object or category types)
 7 categorical_vars = df.select_dtypes(include=['object', 'category']).columns
 8
 9 n_cols = 3
10 n_rows = int(np.ceil(len(categorical_vars) / n_cols))
11 fig, axes = plt.subplots(n_rows, n_cols, figsize=(20, n_rows * 4))
12
13 # Flatten axes array for easy iteration
14 axes = axes.flatten()
15
16 for i, var in enumerate(categorical_vars):
17     sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
18     axes[i].set_title(f'Distribution of {var}', fontsize=12)
19     axes[i].set_xlabel(var, fontsize=10)
20     axes[i].set_ylabel('Count', fontsize=10)
21     axes[i].tick_params(axis='x', rotation=45)
22
23 # Remove any empty subplots if the number of variables is not a multiple of n_cols
24 for j in range(i + 1, len(axes)):
25     fig.delaxes(axes[j])
26
27 plt.tight_layout()
28 plt.show()
```

```
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
<ipython-input-11-2b64ad7672ec>:17: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.countplot(data=df, x=var, ax=axes[i], palette='viridis')
```

## Data Insights

### 1. Gender (RIAGENDR)

#### Categories

- 1 (Male)
- 2 (Female)

**Insight**:
The distribution shows a slightly larger proportion of females (category 2) than males (category 1). This indicates that the dataset is close to balanced by gender, but with a modest female majority.

---

### 2. Race/Ethnicity (RIDRETH1)

#### Categories (Typically coded as)

- Mexican American
- Other Hispanic
- Non-Hispanic White
- Non-Hispanic Black

- Other Race (Including Multi-Racial)

**Insight**:
The largest group appears to be category 2 ("Other Hispanic"), followed by category 3 ("Non-Hispanic White"). Categories 1, 4, and 5 are somewhat smaller but still sizable. This distribution may reflect the specific subset of NHANES participants included in the dataset (or a particular oversampling strategy). Understanding these proportions is crucial when analyzing subgroup trends or health disparities.

## 3. Education Level (DMDEDUC2)

### Categories (Adults 20+)

- less than 9th grade
- 9–11th grade (includes 12th grade with no diploma)
- High school graduate/GED
- Some college or AA degree
- College graduate or above

**Insight**:
"Some college or AA degree" (category 4) has the highest frequency. "High school graduate/GED" (category 3) is next in size, followed by "College graduate or above" (category 5). The smallest categories are those with less than high school education (1 and 2). This suggests that the majority of participants have at least some college experience.

## 4. Marital Status (DMDMARTL)

### Typical Categories

- Married
- Widowed
- Divorced
- Separated
- Never married
- Living with partner

**Insight**:
Married individuals (category 1) form the largest group, closely followed by those who have never married (category 5). Living with a partner (6) is also notable, while widowed (2), divorced (3), and separated (4) are smaller in comparison. This distribution provides a snapshot of social/familial structures that can influence health behaviors and outcomes.

## 5. Country of Birth (DMDBORN4)

### Typical Categories (for NHANES)

- Born in 50 US states or Washington, DC
- Born in a US territory (e.g., Puerto Rico)
- Born in Mexico
- Born outside the US (other than Mexico)

**Insight**:
The largest bar is category 1 (US-born), indicating most participants are native-born. The next largest category is 2 (US territory) or possibly 3 (Mexico), with category 4 also present but smaller. This mix highlights a level of diversity in nativity.
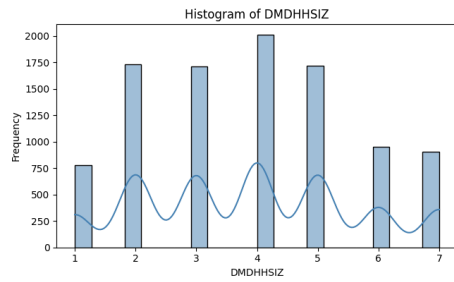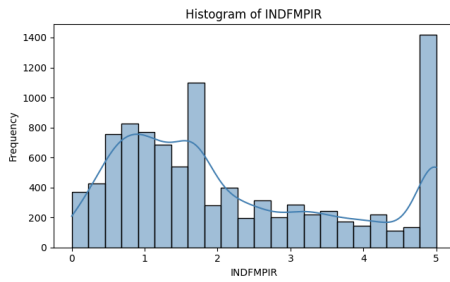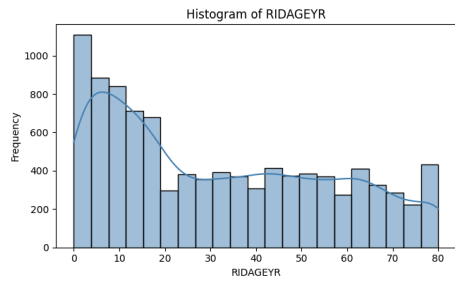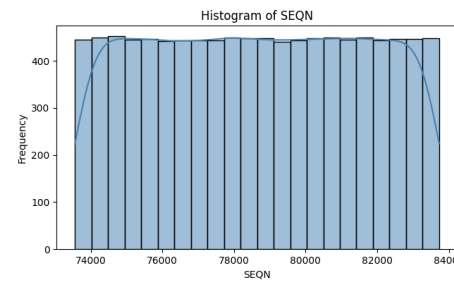
## Overall Observations and Considerations

- **Slight Female Majority**:
  The dataset is almost balanced by gender but leans slightly female. Analyses may need to account for potential gender-related differences in health outcomes.

- **Racial/Ethnic Composition**:
  The distribution shows a strong representation of Hispanic subgroups, with "Other Hispanic" as the largest category. Any subgroup analyses should consider potential oversampling or weighting factors if generalizing to the US population.

- **Educational Attainment**:
  A majority have at least some college education, suggesting relatively high educational attainment. Researchers may want to assess how education level correlates with diet, lifestyle, or mental health outcomes.

- **Marital Status**:
  "Married" and "Never married" are the largest groups. Marital status often correlates with social support and health behaviors, so it may be an important covariate in mental health or chronic disease studies.

- **Nativity**:

  Most participants are US-born, though there is a sizable minority born in other regions. This can be relevant for culturally tailored interventions, immigration-related stress factors, or acculturation studies.

```python
1 # ----------------------------
2 # 3. Continuous Variables Visualization
3 # ----------------------------
4 # Define a list of continuous variables (adjust list based on your dataset)
5 continuous_vars = df.select_dtypes(include=['int64', 'float64']).columns
6
7 n_cols = 3
8 n_rows = int(np.ceil(len(continuous_vars) / n_cols))
9
10 # Histogram with KDE for continuous variables
11 fig, axes = plt.subplots(n_rows, n_cols, figsize=(20, n_rows * 4))
12 axes = axes.flatten()
13
14 for i, var in enumerate(continuous_vars):
15     sns.histplot(data=df, x=var, ax=axes[i], kde=True, color='steelblue')
16     axes[i].set_title(f'Histogram of {var}', fontsize=12)
17     axes[i].set_xlabel(var, fontsize=10)
18     axes[i].set_ylabel('Frequency', fontsize=10)
19
20 # Remove any extra subplots
21 for j in range(i + 1, len(axes)):
22     fig.delaxes(axes[j])
23
24 plt.tight_layout()
25 plt.show()
```

- **Nativity**:

  Most participants are US-born, though there is a sizable minority born in other regions. This can be relevant for culturally tailored interventions, immigration-related stress factors, or acculturation studies.

```python
1 # ----------------------------
```

Histogram of DR1TKCAL · Histogram of DR1TCARB · Histogram of DR1TPROT · Histogram of DR1TTFAT · Histogram of DR1TSFAT · Histogram of DR1TFIBE · Histogram of DR1TSUGR · Histogram of DR1TVD · Histogram of DR1TFOLA · Histogram of DR1TVB12 · Histogram of DR1TIRON · Histogram of DR1TMAGN · Histogram of DR1TCAFF · Histogram of DR1TALCO · Histogram of SEQN · Histogram of RIDAGEYR · Histogram of INDFMPIR · Histogram of DMDHHSIZ

# Dietary Data Analysis

## 1. Total Energy & Macronutrients

**DR1TKCAL (Total Energy Intake)**

- **Shape:** Right-skewed distribution with a concentration around moderate calorie intakes (e.g., ~1,500–2,500 kcal) and a long tail for higher intakes.
- **Insight:** This is typical of dietary data—most individuals cluster around a moderate range, with a small subset reporting very high intakes.

**DR1TCARB (Total Carbohydrate), DR1TPROT (Total Protein), DR1TTFAT (Total Fat), DR1TSFAT (Saturated Fat)**

- **Shape:** All show right-skewed distributions. Most people consume moderate amounts of each nutrient, but a subset has high intakes.
- **Insight:** The skew indicates potential outliers (heavy consumers), which may warrant further inspection or transformation (e.g., log transform) for modeling.

## 2. Fiber and Sugars

**DR1TFIBE (Dietary Fiber), DR1TSUGR (Sugars)**

- **Fiber:** Typically right-skewed, with a peak at lower intake and fewer individuals consuming very high fiber.
- **Sugars:** Also right-skewed, but often with a wider spread; some individuals have very high sugar intakes.
- **Insight:** These distributions reflect varying dietary patterns; high sugar intake may be a point of interest for metabolic or mental health studies.

## 3. Key Micronutrients & Bioactives

**DR1TVD (Vitamin D), DR1TFOLA (Folate), DR1TVB12 (Vitamin B12), DR1TIRON (Iron), DR1TMAGN (Magnesium)**

- **Shape:** All tend to be heavily right-skewed. Many participants cluster at lower intakes, with a minority showing very high values.
- **Insight:** High skew is typical for nutrient data, reflecting differences in supplement use or specific dietary patterns. Researchers should consider whether to treat extreme values as outliers or valid high-consumption cases.

**DR1TCAFF (Caffeine)**

- **Shape:** Wide range from near-zero to high levels of caffeine consumption.
- **Insight:** Caffeine intake can vary drastically based on coffee, tea, soda, or energy drink habits. Extreme values could be relevant for certain health outcomes (e.g., sleep, cardiovascular health).

**DR1TALCO (Alcohol)**

- **Shape:** Most values near zero, with a subset reporting moderate to high alcohol consumption.
- **Insight:** This is typical in population data: many non-drinkers or low-level drinkers, plus a tail of heavier alcohol intake.

## 4. Demographics & Socioeconomics

**RIDAGEYR (Age)**

- **Shape:** Often slightly right-skewed or nearly uniform across adult ages in NHANES-type data, but the histogram suggests a peak in mid-adult ranges and fewer older participants.
- **Insight:** Understanding age distribution is crucial for controlling for age effects in dietary or mental health analyses.

**INDFMPIR (Family Income-to-Poverty Ratio)**

- **Shape:** Often right-skewed, with many participants clustered around lower to moderate income ratios, and a smaller group at higher ratios.
- **Insight:** This variable is critical for socioeconomic status analysis; it can correlate with dietary quality and health outcomes.

**DMDHHSIZ (Household Size)**

- **Shape:** Discrete distribution (1, 2, 3, etc.). Typically peaks around 2–4.
- **Insight:** Larger households could affect resource allocation and dietary behaviors. The histogram may show a strong peak at 2 or 3 members.

## 5. Derived Macronutrient Percentage Metrics

**protein_energy_pct, carb_energy_pct, fat_energy_pct**

- **Shape:** These often cluster around typical dietary proportions (e.g., ~10–20% protein, ~40–60% carbohydrates, ~20–40% fat).
- **Insight:** Distributions may appear more "normal" than raw intake data. Outliers (e.g., extremely high fat-energy percentage) can flag unusual diets.

## Key Observations & Recommendations

### Right-Skew & Potential Outliers:

- Most intake variables show right-skewed distributions. This is common in dietary data but may require transformations (e.g., log) or robust methods to address outliers in statistical models.

### Zero-Inflation:

- Alcohol, caffeine, and certain micronutrients have a large number of zeros (or near-zero) values. Zero-inflated or two-part models might be appropriate in advanced analyses.

### Wide Ranges:

- Caffeine, sugars, and total energy intake exhibit broad distributions, indicating diverse dietary habits. Stratified analyses could reveal subgroups with distinct behaviors or health outcomes.

### Socioeconomic & Demographic Variation:

- Age, family income, and household size vary considerably. These factors often interact with dietary behaviors and should be considered confounders or covariates in regression models.

### Derived Percentages:

- The macronutrient percentage distributions (protein, carb, fat) appear less skewed and can be more interpretable in certain contexts (e.g., recommended dietary guidelines). However, they still require checking for extremes or data quality issues.

### Data Quality & Validity Checks:

- Extremely high or low intakes may be due to reporting errors, so further validation (e.g., 24-hour recall consistency, logic checks) is advisable.

## ⌄ Bivariate Analysis

### Relationship between continuous variables and Target variable

```python
import seaborn as sns
import matplotlib.pyplot as plt

# Define your continuous variables
continuous_vars = [
    'DR1TKCAL', 'DR1TCARB', 'DR1TPROT', 'DR1TTFAT', 'DR1TSFAT',
    'DR1TFIBE', 'DR1TSUGR', 'DR1TVD', 'DR1TFOLA', 'DR1TVB12',
    'DR1TIRON', 'DR1TMAGN', 'DR1TCAFF', 'DR1TALCO', 'RIDAGEYR',
    'INDFMPIR', 'DMDHHSIZ'
]

# Number of rows and columns for subplots (adjust as needed)
n_rows = 5
n_cols = 4

fig, axes = plt.subplots(n_rows, n_cols, figsize=(20, n_rows * 4))
axes = axes.flatten()

# Loop through each continuous variable and create a box plot against the binary target
for i, var in enumerate(continuous_vars):
    sns.boxplot(
        data=df,
        x='mental_health_binary',  # Replace 'y' with the target variable in your dataset
        y=var,
        ax=axes[i],
        palette='viridis'
    )
    axes[i].set_title(f'{var} by Mental Health Status', fontsize=12)
    axes[i].set_xlabel('Mental Health Binary', fontsize=10)
    axes[i].set_ylabel(var, fontsize=10)

# Remove any unused subplots if the number of continuous variables is fewer than n_rows*n_cols
for j in range(i + 1, len(axes)):
    fig.delaxes(axes[j])

plt.tight_layout()
plt.show()
```

```
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
<ipython-input-13-45329984dcc1>:21: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=

  sns.boxplot(
```

DR1TSFAT by Mental Health Status | DR1TFIBE by Mental Health Status | DR1TSUGR by Mental Health Status | DR1TVD by Mental Health Status

DR1TFOLA by Mental Health Status | DR1TVB12 by Mental Health Status | DR1TIRON by Mental Health Status | DR1TMAGN by Mental Health Status

DR1TCAFF by Mental Health Status | DR1TALCO by Mental Health Status | RIDAGEYR by Mental Health Status | INDFMPIR by Mental Health Status

DMDHHSIZ by Mental Health Status

# Comparison of Nutrient Intake and Demographics with Mental Health

## 1. Overall Impressions

- **Similar Medians, Wide Ranges:** Most nutrients and demographic variables show comparable median values between the two groups. This suggests that, at least visually, there are no dramatic mean/median differences.
- **High Variability:** Many box plots feature wide interquartile ranges and numerous outliers, which is common in self-reported dietary data (e.g., total energy, caffeine, and alcohol).

## 2. Nutrient Intake Variables

### Total Energy (DR1TKCAL)

- Both groups display a similar median energy intake, with a skew toward higher outliers.
- This implies that overall caloric consumption does not differ drastically between individuals with and without a mental health condition.

### Macronutrients (DR1TCARB, DR1TPROT, DR1TTFAT, DR1TSFAT)

- Carbohydrate, protein, total fat, and saturated fat intakes show largely overlapping distributions.
- Although there are outliers in each group, there is no clear visual shift in medians or IQRs to indicate a strong association with mental health status.

### Fiber and Sugars (DR1TFIBE, DR1TSUGR)

- **Fiber (DR1TFIBE):** Also appears quite similar across both groups.
- **Sugars (DR1TSUGR):** Might show a slightly higher median for the mental health group (1), but the difference is small, and the spread of values is large.

### Micronutrients (DR1TVD, DR1TFOLA, DR1TVB12, DR1TIRON, DR1TMAGN)

- Vitamin D, folate, vitamin B12, iron, and magnesium are all heavily right-skewed, with numerous outliers in both groups.
- The medians look nearly identical, indicating no major difference in typical micronutrient intake between mental health groups.

### Caffeine (DR1TCAFF) & Alcohol (DR1TALCO)

- **Caffeine:** Both groups have a long-tailed distribution with many outliers, reflecting large variability in coffee, tea, and soda consumption.
- **Alcohol:** Most values are near zero, with a smaller subset reporting moderate to high intakes. The box plots do not suggest a strong difference between groups, though the mental health group appears to have slightly more outliers.

## 3. Demographic & Socioeconomic Variables

### Age (RIDAGEYR)

- The box plots show a broad age range, with both groups having a similar median.
- There is no clear indication that one group is significantly older or younger based on this visualization alone.

### Income-to-Poverty Ratio (INDFMPIR)

- Both distributions appear right-skewed. The mental health group might have a slightly lower median ratio, though the difference is modest.
- Large variability is evident in both groups, suggesting a wide range of socioeconomic statuses.

### Household Size (DMDHHSIZ)

- The median household size is around 2 or 3 in both groups.
- There is some spread, but again no stark difference is immediately apparent.

## 4. Key Observations & Considerations

### No Dramatic Visual Differences:

- From these box plots alone, none of the variables stand out as having a large shift in distribution between the two mental health groups.

### High Outlier Prevalence:

- Dietary data frequently exhibit outliers (e.g., very high sugar or caffeine intakes). Analysts should consider whether these outliers are valid or potential data-entry/reporting errors.

### Potential Minor Trends:

- Sugars and Caffeine may show slightly higher medians among those with mental health conditions, but the differences appear small.

- Income-to-Poverty Ratio might be slightly lower for those with mental health conditions, although again the visual difference is not large.

### Need for Statistical Testing:

- Box plots provide a quick visual check but do not confirm statistical significance. Formal hypothesis testing (e.g., t-tests, non-parametric tests) or modeling (e.g., regression) is necessary to determine whether observed differences are meaningful.

### Multivariable Context:

- Mental health status is influenced by numerous factors (e.g., demographics, genetics, lifestyle).
- Consider adjusting for confounding variables like age, income, and educational level when examining any association between nutrient intake and mental health.
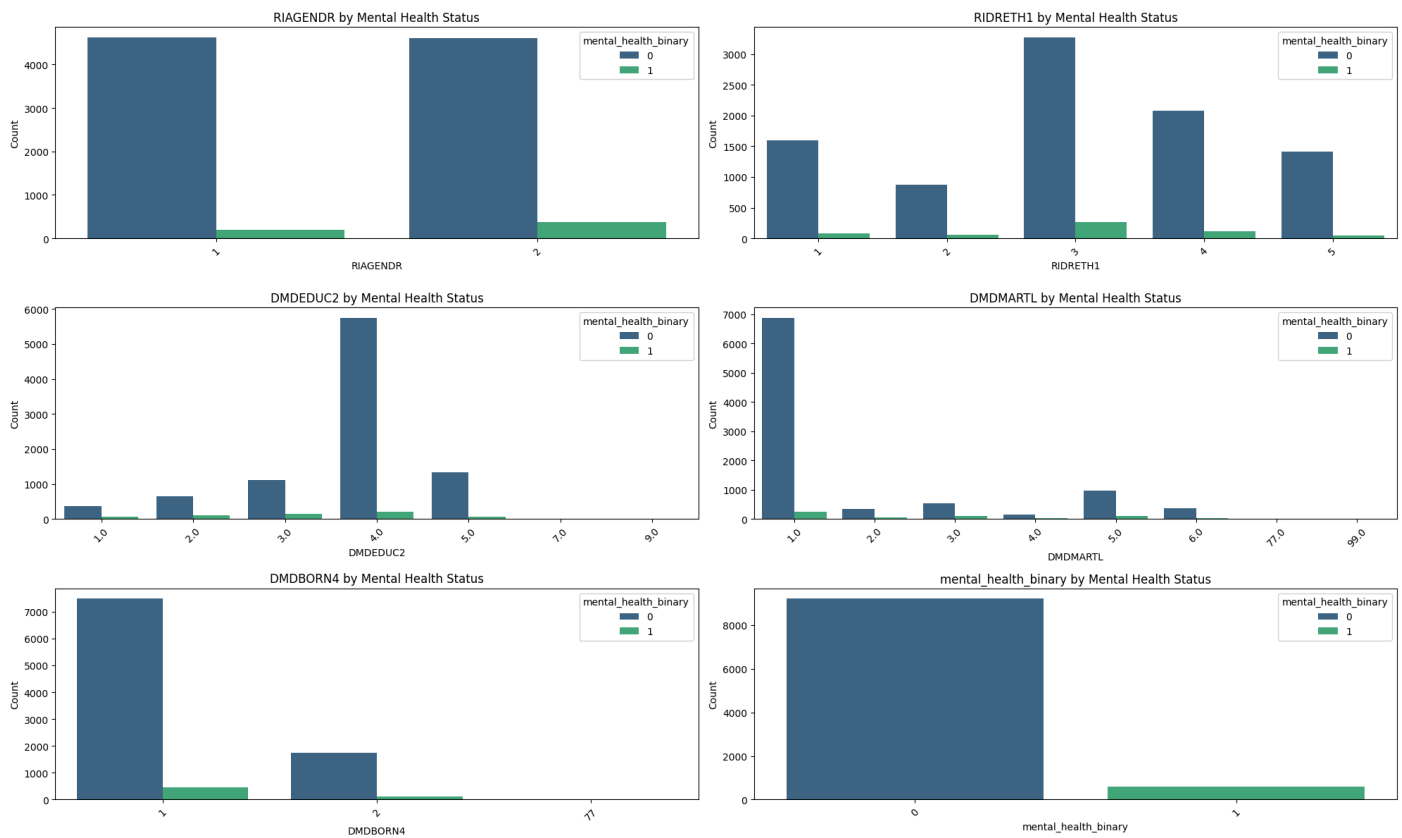
## ∨ Bivariate Analysis

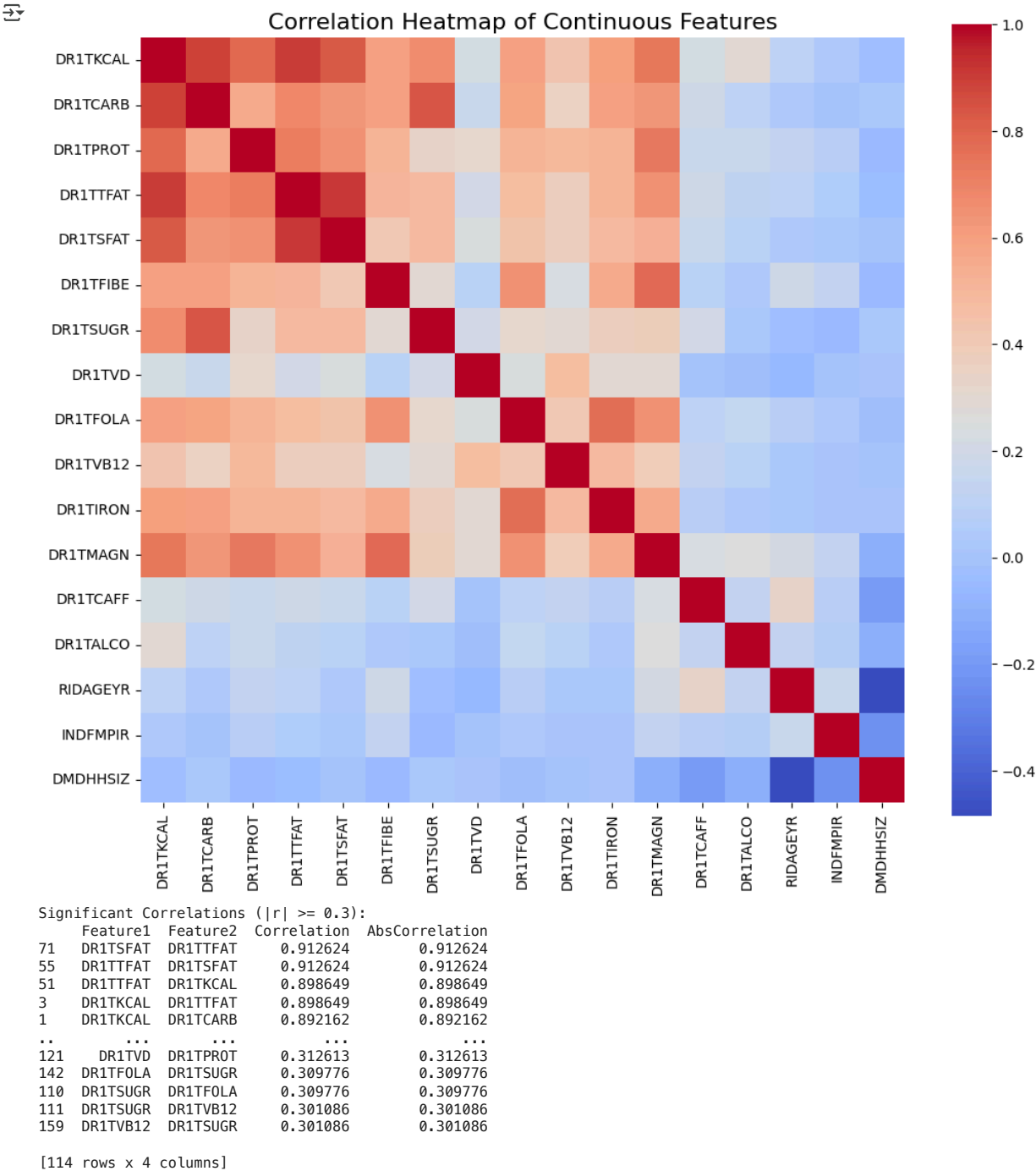### Relationship between categorical variables and Target variable

```
1  import seaborn as sns
2  import matplotlib.pyplot as plt
3
4
5  # Define the number of rows and columns for your subplots
6  n_rows = 3
7  n_cols = 2
8
9  fig, axes = plt.subplots(n_rows, n_cols, figsize=(20, n_rows * 4))
10 axes = axes.flatten()
11
12 for i, var in enumerate(categorical_vars):
13     sns.countplot(
14         data=df,
15         x=var,
16         hue='mental_health_binary',  # Replace 'y' with the target variable
17         ax=axes[i],
18         palette='viridis'
19     )
20     axes[i].set_title(f'{var} by Mental Health Status', fontsize=12)
21     axes[i].set_xlabel(var, fontsize=10)
22     axes[i].set_ylabel('Count', fontsize=10)
23     axes[i].tick_params(axis='x', rotation=45)
24
25 # Remove any extra subplots if you have fewer categorical vars than n_rows*n_cols
26 for j in range(i + 1, len(axes)):
27     fig.delaxes(axes[j])
28
29 plt.tight_layout()
30 plt.show()
31
```

RIAGENDR by Mental Health Status

RIDRETH1 by Mental Health Status

DMDEDUC2 by Mental Health Status

DMDMARTL by Mental Health Status

DMDBORN4 by Mental Health Status

mental_health_binary by Mental Health Status

```
 1 continuous_features = [
 2     'DR1TKCAL', 'DR1TCARB', 'DR1TPROT', 'DR1TTFAT', 'DR1TSFAT',
 3     'DR1TFIBE', 'DR1TSUGR', 'DR1TVD', 'DR1TFOLA', 'DR1TVB12',
 4     'DR1TIRON', 'DR1TMAGN', 'DR1TCAFF', 'DR1TALCO', 'RIDAGEYR',
 5     'INDFMPIR', 'DMDHHSIZ'
 6 ]
 7
 8 # Subset the DataFrame to include only these columns
 9 df_continuous = df[continuous_features].copy()
10
11 # ----------------------------------------------------
12 # Compute the Correlation Matrix
13 # ----------------------------------------------------
14 corr_matrix = df_continuous.corr()
15
16 # ----------------------------------------------------
17 #  Create a Correlation Heatmap
18 # ----------------------------------------------------
19 plt.figure(figsize=(12, 10))
20 sns.heatmap(corr_matrix, annot=False, cmap='coolwarm', square=True)
21 plt.title('Correlation Heatmap of Continuous Features', fontsize=16)
22 plt.show()
23
24 # ----------------------------------------------------
25 # Summarize Significant Correlations
26 # ----------------------------------------------------
27 # Flatten the correlation matrix and sort by absolute value
28 corr_pairs = corr_matrix.unstack().reset_index()
29 corr_pairs.columns = ['Feature1', 'Feature2', 'Correlation']
30
31 # Remove self-correlations (where Feature1 == Feature2)
32 corr_pairs = corr_pairs[corr_pairs['Feature1'] != corr_pairs['Feature2']]
33
34 # Convert to absolute correlation values for filtering and sort descending
35 corr_pairs['AbsCorrelation'] = corr_pairs['Correlation'].abs()
36 corr_pairs = corr_pairs.sort_values('AbsCorrelation', ascending=False)
37
38 # Define a threshold for "significant" correlation (e.g., |r| >= 0.3)
39 threshold = 0.3
40 significant_corr = corr_pairs[corr_pairs['AbsCorrelation'] >= threshold]
41
42 print("Significant Correlations (|r| >= 0.3):")
43 print(significant_corr)
```

## Correlation Heatmap of Continuous Features



```
Significant Correlations (|r| >= 0.3):
      Feature1  Feature2  Correlation  AbsCorrelation
71    DR1TSFAT  DR1TTFAT     0.912624        0.912624
55    DR1TTFAT  DR1TSFAT     0.912624        0.912624
51    DR1TTFAT  DR1TKCAL     0.898649        0.898649
3     DR1TKCAL  DR1TTFAT     0.898649        0.898649
1     DR1TKCAL  DR1TCARB     0.892162        0.892162
..         ...       ...          ...             ...
121    DR1TVD   DR1TPROT     0.312613        0.312613
142   DR1TFOLA  DR1TSUGR     0.309776        0.309776
110   DR1TSUGR  DR1TFOLA     0.309776        0.309776
111   DR1TSUGR  DR1TVB12     0.301086        0.301086
159   DR1TVB12  DR1TSUGR     0.301086        0.301086

[114 rows x 4 columns]
```

## ∨ Key Observations from the Heatmap

### Total Energy vs. Macronutrients

- **DR1TKCAL (Total Energy Intake)** shows strong positive correlations with the absolute amounts of carbohydrate (**DR1TCARB**), protein (**DR1TPROT**), total fat (**DR1TTFAT**), saturated fat (**DR1TSFAT**), and sugars (**DR1TSUGR**).
- **Insights:** The more total calories someone consumes, the more of each macronutrient (carbs, protein, fat) they tend to eat overall.

### Saturated Fat and Total Fat

- **DR1TSFAT (Saturated Fat)** is strongly correlated with **DR1TTFAT (Total Fat)**.
- **Insights:** People who eat more total fat also typically eat more saturated fat, which is a subset of total fat.

### Carbohydrates and Sugars

- **DR1TCARB (Carbohydrates)** is strongly positively correlated with **DR1TSUGR (Sugars)** and moderately with **DR1TFIBE (Fiber)**.
- **Insights:** Individuals who consume more total carbohydrates also tend to consume more sugar and more fiber.

## Micronutrients & Others

- Micronutrients like **Vitamin D (DR1TVD)**, **Folate (DR1TFOLA)**, **Vitamin B12 (DR1TVB12)**, **Iron (DR1TIRON)**, and **Magnesium (DR1TMAGN)** show some moderate correlations with total calorie intake or each other, but none appear extremely strong in the heatmap.
- **Caffeine (DR1TCAFF)** and **Alcohol (DR1TALCO)** do not show particularly strong correlations with other nutrients or demographic factors, suggesting these intakes may be more independent lifestyle choices.

## Demographics & Socioeconomics

- **RIDAGEYR (Age)**, **INDFMPIR (Income-to-Poverty Ratio)**, and **DMDHHSIZ (Household Size)** do not appear to have strong correlations with most dietary variables.
- **Insights:** Being older or younger, or having a larger or smaller household, does not necessarily mean you consume drastically different amounts of these nutrients. Family income ratio also does not show a very strong linear relationship with these intake measures in this sample.

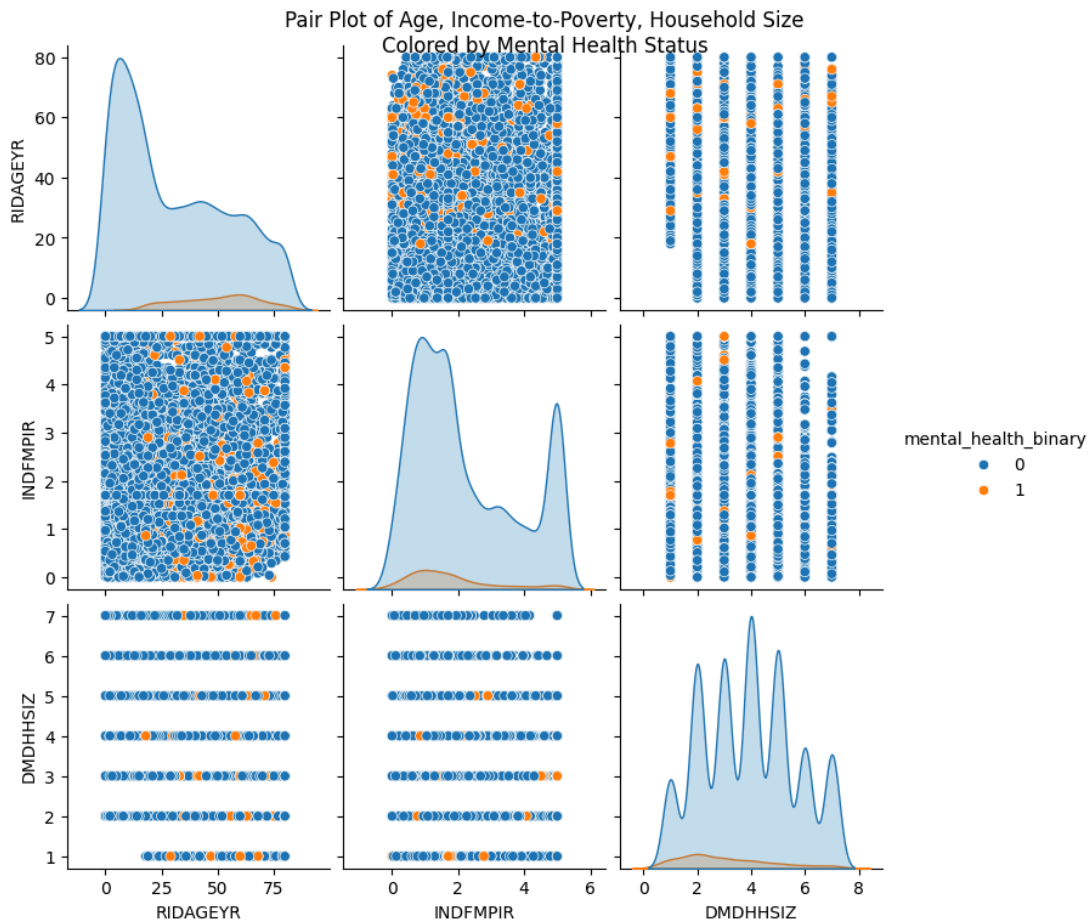## Significant Correlations Insights

Here are a few of the more notable correlations, described simply:

- **Total Energy & Carbohydrates:** "People who eat more calories overall also tend to eat more total carbohydrates."
- **Saturated Fat & Total Fat:** "When someone's total fat intake goes up, their saturated fat intake usually goes up too."
- **Carbohydrates & Sugars:** "If you eat a lot of carbohydrates, you usually end up eating more sugar as well."
- **Micronutrients:** "Vitamins and minerals (like Vitamin D, Iron, Folate) have some moderate relationships with overall calorie intake, but none are extremely high."

Double-click (or enter) to edit

```
 1 import seaborn as sns
 2 import matplotlib.pyplot as plt
 3
 4 # Example subset of columns to explore:
 5 #   - RIDAGEYR (Age)
 6 #   - INDFMPIR (Income-to-Poverty Ratio)
 7 #   - DMDHHSIZ (Household Size)
 8 #   - mental_health_binary (Binary Target)
 9 # Note: Ensure mental_health_binary is a categorical dtype.
10 df['mental_health_binary'] = df['mental_health_binary'].astype('category')
11
12 pairplot_vars = ['RIDAGEYR', 'INDFMPIR', 'DMDHHSIZ']
13
14 sns.pairplot(
15     data=df,
16     vars=pairplot_vars,
17     hue='mental_health_binary',   # Color-code by mental health status
18     diag_kind='kde'               # Use KDE (Kernel Density Estimate) on the diagonal
19 )
20
21 plt.suptitle("Pair Plot of Age, Income-to-Poverty, Household Size\nColored by Mental Health Status", y=1.02)
22 plt.show()
23
```
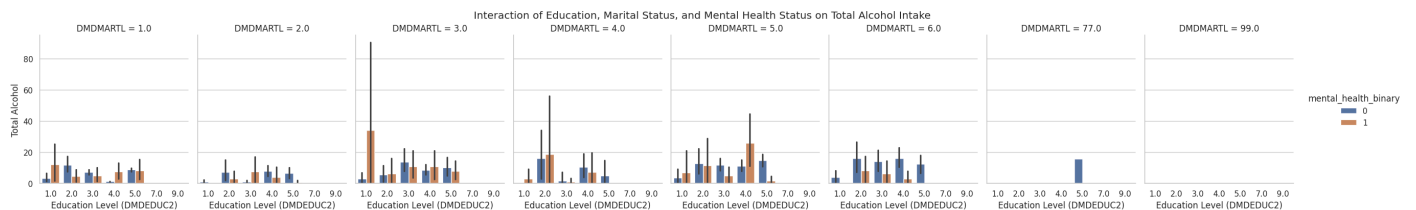
Pair Plot of Age, Income-to-Poverty, Household Size
Colored by Mental Health Status

Double-click (or enter) to edit

```
1  # Example catplot showing how total energy intake (DR1TKCAL) varies
2  # by education level (DMDEDUC2) and marital status (DMDMARTL),
3  # stratified by mental health status (mental_health_binary).
4
5  sns.set_theme(style="whitegrid")
6
7  g = sns.catplot(
8      data=df,
9      x='DMDEDUC2',              # Education category on the x-axis
10     y='DR1TALCO',              # Numeric outcome (Total Alcohol intake)
11     hue='mental_health_binary',  # Separate colors by mental health status
12     col='DMDMARTL',            # Facet (column) by marital status
13     kind='bar',                # Bar plot (can also use 'point', 'box', etc.)
14     height=4, aspect=0.8
15  )
16
17  g.fig.suptitle("Interaction of Education, Marital Status, and Mental Health Status on Total Alcohol Intake", y=1.02)
18  g.set_axis_labels("Education Level (DMDEDUC2)", "Total Alcohol")
19  plt.show()
20
```



## Alcohol Consumption Analysis

## 1. Overall Alcohol Consumption

## Mostly Low or Near Zero

- In most columns, the bars are close to zero, suggesting that many individuals in these groups do not consume much alcohol on a daily basis (or do not drink at all).

## Occasional Higher Averages

- A few groups have taller bars, which could indicate smaller sample sizes (fewer people in that category) or a handful of individuals reporting higher alcohol consumption.

## 2. Education Level (DMDEDUC2)

Typical NHANES-like codes might be:

- Less than 9th grade
- 9–11th grade (no diploma)
- High school graduate/GED
- Some college or AA degree
- College graduate or above
- (Other codes, like 9, might be unknown/missing.)

### No Clear Trend by Education

- There is no obvious pattern where more or less schooling consistently relates to higher or lower alcohol intake. Averages remain low in most education groups, with a few spikes that appear sporadic.

## 3. Marital Status (DMDMARTL)

Typical codes could be:

- Married
- Widowed
- Divorced
- Separated
- Never Married
- Living with partner
- 77/99. Refused / Don't Know

### Similar Low Levels

- Whether married, divorced, never married, or living with a partner, most bars indicate relatively low daily alcohol consumption on average.

### Small Groups

- Categories like 77 or 99 (Refused/Don't Know) might show unusual heights because they contain fewer participants.

## 4. Mental Health Status Differences

### Orange vs. Blue Bars

- The difference between mental health status 1 (orange) and status 0 (blue) is not dramatic in most columns. Some are slightly higher or lower, but there's no uniform pattern across all marital and education levels.

### Minor Variations

- Where bars differ, the gap is generally small, suggesting that having a mental health condition does not strongly drive daily alcohol intake in these groups.

## 5. Summary

- **Low Alcohol Intake Overall:** Most people, regardless of education level, marital status, or mental health, report little to no daily alcohol consumption.
- **Occasional Higher Spikes:** Certain categories spike higher, but this is likely due to fewer participants or a few individuals who drink more.
- **No Major Group Gaps:** Education and marital status do not show a clear, consistent effect on drinking habits. Similarly, those with and without a mental health condition often have comparable alcohol intake.
- **Further Analysis:** If you want to confirm whether small differences are meaningful, you'd need additional statistical tests or more in-depth study. However, at a glance, this chart indicates that most people in the dataset do not consume large amounts of alcohol daily, and demographic factors alone do not appear to dramatically change that behavior.

```
1  import pandas as pd
2  import numpy as np
3  from scipy.stats import ttest_ind, mannwhitneyu
```

```
 4
 5  # Example DataFrame: df
 6  # Numeric variable to compare (e.g., total energy intake)
 7  numeric_var = 'DR1TKCAL'          # Adjust to your numeric column
 8  # Binary grouping variable (e.g., mental health status)
 9  group_var = 'mental_health_binary' # Adjust if your target column has a different name
10
11  # Separate the data into two groups based on the binary variable
12  group0 = df[df[group_var] == 0][numeric_var].dropna()
13  group1 = df[df[group_var] == 1][numeric_var].dropna()
14
15  ##########################################
16  # 1. Independent Samples T-test
17  ##########################################
18  # Assumes the data is (approximately) normally distributed in each group
19  # and the observations are independent.
20
21  t_stat, p_val = ttest_ind(group0, group1, equal_var=False)
22  print(f"T-test for {numeric_var} by {group_var}:")
23  print(f"t-statistic = {t_stat:.3f}, p-value = {p_val:.3f}")
24
25  ##########################################
26  # 2. Mann-Whitney U Test
27  ##########################################
28  # A non-parametric test that does not assume normal distribution.
29
30  u_stat, p_val_mw = mannwhitneyu(group0, group1, alternative='two-sided')
31  print(f"\nMann-Whitney U test for {numeric_var} by {group_var}:")
32  print(f"U-statistic = {u_stat:.3f}, p-value = {p_val_mw:.3f}")
33
```

```
⤓  T-test for DR1TKCAL by mental_health_binary:
     t-statistic = -1.937, p-value = 0.053

     Mann-Whitney U test for DR1TKCAL by mental_health_binary:
     U-statistic = 2643701.500, p-value = 0.403
```

## ⌄ Statistical Test Results

Based on the results, here's what we can conclude:

T-test Result:

- **t-statistic:** -1.937
- **p-value:** 0.053
- The p-value is just above the conventional **0.05 threshold**. This means that the difference in total energy intake (**DR1TKCAL**) between the two mental health groups is almost significant, but not quite.
- It suggests a borderline or marginal difference—if we had a slightly larger sample size or less variability, we might see a significant effect.

Mann-Whitney U Test Result:

- **U-statistic:** 2643701.500
- **p-value:** 0.403
- This non-parametric test shows a much higher p-value, which strongly indicates that there is no statistically significant difference between the groups.

## ⌄ Applying Transformations and Deriving new features

```
 1  ##########################################
 2  # 1. Derived Macronutrient Energy Contributions
 3  ##########################################
 4  # Note: Protein and Carbohydrates provide 4 kcal/g; Fat provides 9 kcal/g.
 5  df['protein_energy'] = df['DR1TPROT'] * 4
 6  df['carb_energy'] = df['DR1TCARB'] * 4
 7  df['fat_energy'] = df['DR1TTFAT'] * 9
 8
 9  # Calculate the percentage contribution of each macronutrient to total energy.
10  # Use np.where to avoid division by zero.
11  df['protein_pct'] = np.where(df['DR1TKCAL'] > 0, df['protein_energy'] / df['DR1TKCAL'], np.nan)
12  df['carb_pct'] = np.where(df['DR1TKCAL'] > 0, df['carb_energy'] / df['DR1TKCAL'], np.nan)
13  df['fat_pct'] = np.where(df['DR1TKCAL'] > 0, df['fat_energy'] / df['DR1TKCAL'], np.nan)
14
15  ##########################################
16  # 2. Log Transformation for Skewed Variables
17  ##########################################
18  # Many nutrient intake variables are right-skewed.
```

```python
18 # Many nutrient intake variables are right skewed.
19 # Log transformation (using log1p to handle zeros) can help stabilize variance.
20 df['log_DR1TKCAL'] = np.log1p(df['DR1TKCAL'])
21 df['log_DR1TCARB'] = np.log1p(df['DR1TCARB'])
22 df['log_DR1TCAFF'] = np.log1p(df['DR1TCAFF'])
23 # Add additional log transformations for other variables if needed.
24
25 #########################################
26 # 3. Creating Binary Features Based on Quantiles
27 #########################################
28 # Example: Create a binary feature for "high sugar intake" using the 75th percentile.
29 sugar_threshold = df['DR1TSUGR'].quantile(0.75)
30 df['high_sugar'] = np.where(df['DR1TSUGR'] >= sugar_threshold, 1, 0)
31
32 # Similarly, for alcohol intake:
33 alcohol_threshold = df['DR1TALCO'].quantile(0.75)
34 df['high_alcohol'] = np.where(df['DR1TALCO'] >= alcohol_threshold, 1, 0)
35
36 #########################################
37 # 4. Recoding Categorical Variables to Meaningful Labels
38 #########################################
39 # Example for Gender (assuming 1 = Male, 2 = Female)
40 df['gender'] = df['RIAGENDR'].replace({1: 'Male', 2: 'Female'})
41
42 # Example for Education Level (DMDEDUC2)
43 education_map = {
44     1: '<9th grade',
45     2: '9-11th grade',
46     3: 'High school graduate/GED',
47     4: 'Some college/AA degree',
48     5: 'College graduate or above'
49 }
50 df['education'] = df['DMDEDUC2'].map(education_map)
51
52 # Example for Marital Status (DMDMARTL)
53 marital_map = {
54     1: 'Married',
55     2: 'Widowed',
56     3: 'Divorced',
57     4: 'Separated',
58     5: 'Never married',
59     6: 'Living with partner',
60     77: 'Refused',
61     99: "Don't know"
62 }
63 df['marital_status'] = df['DMDMARTL'].map(marital_map)
64
65 #########################################
66 # 5. Creating Age Groups
67 #########################################
68 # Create an age group variable based on RIDAGEYR.
69 # Adjust bins and labels as needed.
70 df['age_group'] = pd.cut(df['RIDAGEYR'],
71                          bins=[0, 18, 35, 50, 65, 100],
72                          labels=['Child', 'Young Adult', 'Adult', 'Mid-Age', 'Senior'],
73                          right=False)
74
75 # Print out the first few rows to verify the new features
76 print("Feature engineering complete. Preview of new features:")
77 print(df[['DR1TKCAL', 'protein_energy', 'carb_energy', 'fat_energy',
78           'protein_pct', 'carb_pct', 'fat_pct', 'log_DR1TKCAL',
79           'high_sugar', 'high_alcohol', 'gender', 'education', 'marital_status', 'age_group']].head())
80
```

```
Feature engineering complete. Preview of new features:
    DR1TKCAL  protein_energy  carb_energy  fat_energy  protein_pct  carb_pct  \
0     1574.0          174.52       958.36      475.29     0.110877  0.608869
1     5062.0         1352.52      1695.12     1118.61     0.267191  0.334872
2     1743.0          258.44       897.56      593.73     0.148273  0.514951
3     1490.0          311.00       651.68      524.43     0.208725  0.437369
4     1421.0          220.96       712.80      498.24     0.155496  0.501619

    fat_pct  log_DR1TKCAL  high_sugar  high_alcohol  gender  \
0  0.301963      7.362011           1             1    Male
1  0.220982      8.529714           0             1    Male
2  0.340637      7.463937           0             1    Male
3  0.351966      7.307202           0             1    Male
4  0.350626      7.259820           0             1  Female

                   education marital_status age_group
0  High school graduate/GED      Separated    Senior
1  High school graduate/GED        Married   Mid-Age
2    Some college/AA degree        Married    Senior
3    Some college/AA degree        Married     Child
4  College graduate or above       Married    Senior
<ipython-input-19-65406a2910fc>:40: FutureWarning: The behavior of Series.replace (and DataFrame.replace) with CategoricalDtype is deprec
  df['gender'] = df['RIAGENDR'].replace({1: 'Male', 2: 'Female'})
```

## Applying StandardScaler & OneHot Encoding

```python
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4
5 # Assume df is already loaded and preprocessed (from previous steps)
6 # Define the list of numeric columns (e.g., nutrient intakes and demographic variables)
7 numeric_cols = [
8     'DR1TKCAL', 'DR1TCARB', 'DR1TPROT', 'DR1TTFAT', 'DR1TSFAT',
9     'DR1TFIBE', 'DR1TSUGR', 'DR1TVD', 'DR1TFOLA', 'DR1TVB12',
10    'DR1TIRON', 'DR1TMAGN', 'DR1TCAFF', 'DR1TALCO', 'RIDAGEYR',
11    'INDFMPIR', 'DMDHHSIZ', 'protein_energy', 'carb_energy', 'fat_energy',
12    'protein_pct', 'carb_pct', 'fat_pct', 'log_DR1TKCAL',
13    'high_sugar', 'high_alcohol'
14 ]
15
16 # Define the list of categorical columns (using recoded labels or original codes)
17 categorical_cols = [
18     'gender',          # Derived from RIAGENDR (e.g., 'Male', 'Female')
19     'education',       # Derived from DMDEDUC2 (e.g., '<9th grade', 'High school graduate/GED', etc.)
20     'marital_status',# Derived from DMDMARTL (e.g., 'Married', 'Divorced', etc.)
21     'age_group',       # Derived from RIDAGEYR (e.g., 'Young Adult', 'Adult', etc.)
22 ]
23
24 # Create a ColumnTransformer that scales numeric features and one-hot encodes categorical features
25 preprocessor = ColumnTransformer(
26     transformers=[
27         ('num', StandardScaler(), numeric_cols),
28         ('cat', OneHotEncoder(drop='first'), categorical_cols)
29     ]
30 )
31
32 # Apply the transformations (fit and transform the data)
33 df_transformed = preprocessor.fit_transform(df)
34
35 # The resulting df_transformed is a NumPy array.
36 print("Transformed Data Shape:", df_transformed.shape)
37
```

```
Transformed Data Shape: (9813, 43)
```

## Modeling

Logistic Regression as Baseline Model

```python
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.pipeline import Pipeline
6 from sklearn.compose import ColumnTransformer
7 from sklearn.preprocessing import StandardScaler, OneHotEncoder
8 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
9
10
11 # ----------------------------
12 # 2. Define Features and Target
13 # ----------------------------
14 numeric_cols = [
15     'DR1TKCAL', 'DR1TCARB', 'DR1TPROT', 'DR1TTFAT', 'DR1TSFAT',
16     'DR1TFIBE', 'DR1TSUGR', 'DR1TVD', 'DR1TFOLA', 'DR1TVB12',
17     'DR1TIRON', 'DR1TMAGN', 'DR1TCAFF', 'DR1TALCO', 'RIDAGEYR',
18     'INDFMPIR', 'DMDHHSIZ', 'protein_energy', 'carb_energy', 'fat_energy',
19     'protein_pct', 'carb_pct', 'fat_pct', 'log_DR1TKCAL',
20     'high_sugar', 'high_alcohol'
21 ]
22
23 # Define the list of categorical columns (using recoded labels or original codes)
24 categorical_cols = [
25     'gender',          # Derived from RIAGENDR (e.g., 'Male', 'Female')
26     'education',       # Derived from DMDEDUC2 (e.g., '<9th grade', 'High school graduate/GED', etc.)
27     'marital_status',# Derived from DMDMARTL (e.g., 'Married', 'Divorced', etc.)
28     'age_group',       # Derived from RIDAGEYR (e.g., 'Young Adult', 'Adult', etc.)
29 ]
30
31 target = 'mental_health_binary'
32
33 # ----------------------------
```

```
34 # 3. Create a Preprocessing Pipeline
35 # ----------------------------
36 preprocessor = ColumnTransformer(
37     transformers=[
38         ('num', StandardScaler(), numeric_cols),
39         ('cat', OneHotEncoder(drop='first'), categorical_cols)
40     ]
41 )
42
43 # ----------------------------
44 # 4. Split the Data into Training and Testing Sets
45 # ----------------------------
46 X = df[numeric_cols + categorical_cols]
47 y = df[target]
48
49 X_train, X_test, y_train, y_test = train_test_split(X, y,
50                                                     test_size=0.3,
51                                                     random_state=42,
52                                                     stratify=y)
53
54 # ----------------------------
55 # 5. Build a Pipeline with Logistic Regression
56 # ----------------------------
57 pipeline = Pipeline(steps=[
58     ('preprocessor', preprocessor),
59     ('classifier', LogisticRegression(solver='liblinear', random_state=42))
60 ])
61
62 # ----------------------------
63 # 6. Train the Model
64 # ----------------------------
65 pipeline.fit(X_train, y_train)
66
67 # ----------------------------
68 # 7. Evaluate the Model
69 # ----------------------------
70 y_pred = pipeline.predict(X_test)
71
72 print("Confusion Matrix:")
73 print(confusion_matrix(y_test, y_pred))
74
```

```
Confusion Matrix:
[[2768    0]
 [ 171    5]]
```

```
 1 import pandas as pd
 2 from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score, accuracy_score
 3
 4 # Obtain predictions and predicted probabilities for the positive class
 5 y_pred = pipeline.predict(X_test)
 6 y_pred_proba = pipeline.predict_proba(X_test)[:, 1]
 7
 8 # Calculate evaluation metrics for the positive class (assumed label=1)
 9 precision = precision_score(y_test, y_pred, pos_label=1)
10 recall = recall_score(y_test, y_pred, pos_label=1)
11 f1 = f1_score(y_test, y_pred, pos_label=1)
12 auc = roc_auc_score(y_test, y_pred_proba)
13 accuracy = accuracy_score(y_test, y_pred)
14
15 # Create a DataFrame to display the results as a table
16 results_df = pd.DataFrame({
17     'Model': ['Logistic Regression'],
18     'Precision': [precision],
19     'Recall': [recall],
20     'F1-Score': [f1],
21     'AUC': [auc],
22     'Accuracy': [accuracy]
23 })
24
25 results_df
26
```

| | Model | Precision | Recall | F1-Score | AUC | Accuracy |
|---|---|---|---|---|---|---|
| 0 | Logistic Regression | 1.0 | 0.028409 | 0.055249 | 0.819662 | 0.941916 |

## ⌄ Evaluation Metrics Analysis

## 1. Valid Interpretation of an Evaluation Metric

## Precision (1.0):

- Precision measures the proportion of instances predicted as positive that are truly positive. A precision of **1.0** means that every case the model flagged as having a mental health condition was indeed a true positive. However, because precision does not account for missed positives (false negatives), this value must be interpreted alongside other metrics.

## Recall (0.028):

- Recall, also known as **sensitivity**, is the proportion of actual positive cases that were correctly identified by the model. In the results, a recall of **0.028** means that the model identified only about **2.8%** of the individuals who actually have a mental health condition. This clearly indicates that, despite the perfect precision, the model is missing a vast majority of positive cases.

## Precision:

- **Rationale:** In healthcare applications, high precision is important when the cost of a false positive is high (e.g., unnecessary stress for a patient or unwarranted follow-up tests). Here, a precision of **1.0** tells us that if the model flags someone as having a mental health condition, we can be very confident that this prediction is correct.

## Recall:

- **Rationale:** Recall is critical in healthcare, especially when the risk of missing a diagnosis is unacceptable. A very low recall (**2.8%** in this case) is concerning because it indicates that the model is not capturing most of the individuals with the condition. This could lead to under-diagnosis and missed opportunities for early intervention.

## F1-Score (0.055):

- **Rationale:** The **F1-Score** combines precision and recall into a single metric using their harmonic mean. It provides a more balanced view of the model performance when dealing with imbalanced datasets. In your case, the low F1-Score reflects the model's poor ability to identify positive cases, despite high precision.

## AUC (0.819662):

- **Rationale:** The **Area Under the ROC Curve (AUC)** measures the model's overall ability to distinguish between classes, regardless of the chosen threshold. An AUC of about **0.82** suggests that, when considering the ranking of cases by predicted probability, the model has a good discriminative ability. This is important for understanding the model's potential, even if the fixed threshold performance (as reflected by recall) is poor.

## Accuracy (0.941916):

- **Rationale: Accuracy** indicates the overall proportion of correct predictions. However, in imbalanced datasets—such as one where the positive class (mental health condition) is rare—a high accuracy can be misleading. The high accuracy here is driven by the model correctly identifying the majority negative class, while missing almost all positive cases.

## Summary

- **Precision:** Tells us that if the model flags someone as having a mental health issue, it's almost certainly right.
- **Recall:** Tells us that the model is missing almost all the people who actually have the issue, capturing only about **3%** of them.
- **F1-Score:** Provides a combined picture of these two metrics, and the low value indicates that the model isn't very useful in practice despite its high precision.
- **AUC:** Suggests that the model can rank individuals by risk fairly well, even though the fixed threshold used in classification leads to low recall.
- **Accuracy:** Appears high because most individuals do not have the condition, but this masks the fact that the model fails to identify the positives.

```
1 Start coding or generate with AI.
```

```
1 Start coding or generate with AI.
```

```
1 Start coding or generate with AI.
```

## ∨ Evaluation Metrics Analysis

### 1. Valid Interpretation of an Evaluation Metric

## Precision (1.0):

- Precision measures the proportion of instances predicted as positive that are truly positive. A precision of **1.0** means that every case the model flagged as having a mental health condition was indeed a true positive. However, because precision does not account for missed positives (false negatives), this value must be interpreted alongside other metrics.