

```
In [17]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
import pandas as pd
```

```
In [18]: import seaborn as sns
```

```
In [19]: # Load the dataset
df = pd.read_csv('depression_data.csv')
df.head()
```

	Name	Age	Marital Status	Education Level	Number of Children	Smoking Status	Physical Activity Level	Employment Status	Income	Alcohol Consumption	Dietary Habits	Sleep Patterns	History of Mental Illness	History of Substance Abuse	Family History of Depression	Chronic Medical Conditions
0	Christine Baker	31	Married	Bachelor's Degree	2	Non-smoker	Active	Unemployed	35355.67	Moderate	Moderate	Fair	Yes	No	No	Yes
1	Zimmerman Lewis	55	Married	High School	1	Non-smoker	Sedentary	Employed	42710.38	High	Unhealthy	Fair	Yes	No	No	Yes
2	Shannon Chanen	78	Widowed	Master's Degree	1	Non-smoker	Sedentary	Employed	125332.79	Low	Unhealthy	Good	No	No	Yes	No
3	Charles Jordan	58	Divorced	Master's Degree	3	Non-smoker	Moderate	Unemployed	9592.78	Moderate	Moderate	Poor	No	No	No	No
4	Michael Rich	18	Single	High School	0	Non-smoker	Sedentary	Unemployed	8595.08	Low	Moderate	Fair	Yes	No	Yes	Yes

```
In [20]: # Information of the dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 43768 entries, 0 to 43767
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Name                  43768 non-null   object
 1   Age                   43768 non-null   int64
 2   Marital Status        43768 non-null   object
 3   Education Level        43768 non-null   object
 4   Number of Children     43768 non-null   int64
 5   Smoking Status         43768 non-null   object
 6   Physical Activity Level 43768 non-null   object
 7   Employment Status      43768 non-null   object
 8   Income                 43768 non-null   float64
 9   Alcohol Consumption    43768 non-null   object
10   Dietary Habits         43768 non-null   object
11   Sleep Patterns         43768 non-null   object
12   History of Substance Abuse 43768 non-null  object
13   Family History of Depression 43768 non-null  object
14   Chronic Medical Conditions 43768 non-null  object
15   Chronic Medical Conditions_Yes 43768 non-null  object
dtypes: float64 1, int64(2), object(13)
memory usage: 30.5+ MB
```

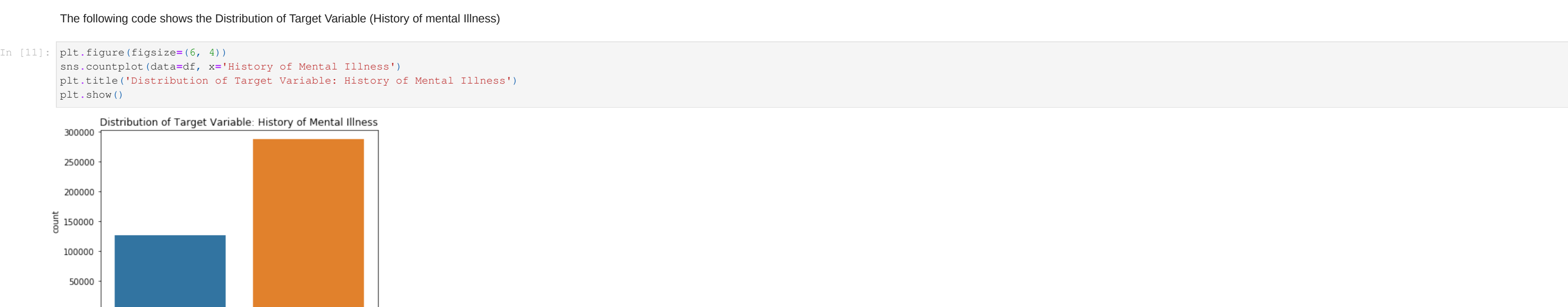
```
In [21]: #Describe the state of the dataset
df.describe()
```

	Age	Number of Children	Income
count	43768.000000	43768.000000	43768.000000
mean	49.000713	1.298972	50661.70971
std	18.158759	1.237054	40624.10066
min	18.000000	0.000000	0.410000
25%	33.000000	0.000000	21001.000000
50%	49.000000	1.000000	67520.150000
75%	65.000000	2.000000	76816.300000
max	80.000000	4.000000	209965.220000

```
In [22]: #Boxplot for missing values
df.isnull().sum()
```

Name	0
Age	0
Marital Status	0
Education Level	0
Number of Children	0
Smoking Status	0
Physical Activity Level	0
Employment Status	0
Income	0
Alcohol Consumption	0
Dietary Habits	0
Sleep Patterns	0
History of Mental Illness	0
History of Substance Abuse	0
Family History of Depression	0
Chronic Medical Conditions	0
dtypes	int64

The following code shows the Distribution of Target Variable (History of mental illness)



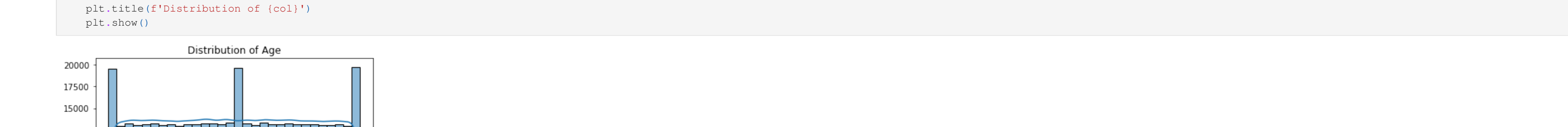
**Observation:** As we can see from the above that there is a clear imbalance with individuals with no history of mental illness compared to do those who do. There is approximately 70% data of individuals with no history of mental illness

**Implication:** This suggests that we could use SMOTE resampling techniques for more accurate predictors

**Distribution of Numerical values**

```
In [10]: numerical_cols = ['Age', 'Number of Children', 'Income']

for col in numerical_cols:
    plt.figure(figsize=(5, 4))
    sns.histplot(df[col], bins=30, kde=True)
    plt.title(f'Distribution of {col}')
    plt.show()
```



**Insight-1:** There is a uniform distribution in the age group which suggests that it is well balanced; The spikes indicate trends within specific age groups related to the events

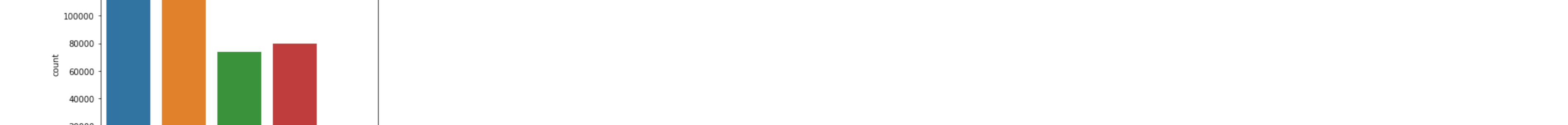
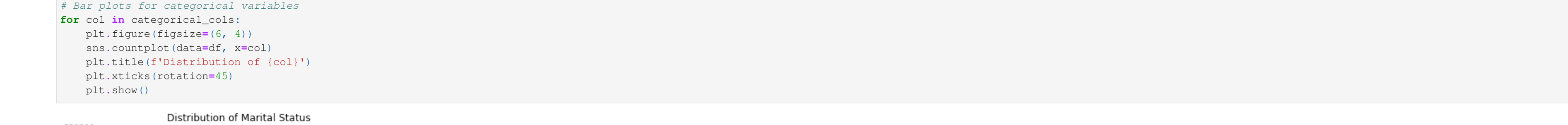
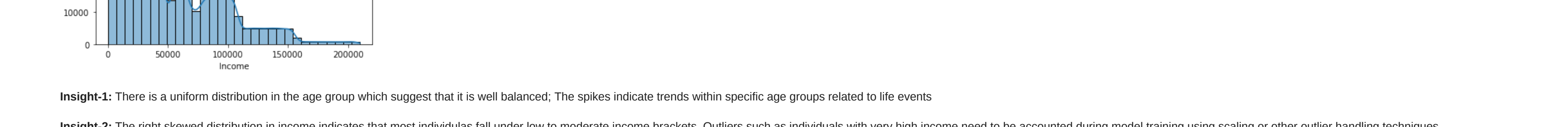
**Insight-2:** The right skewed distribution in income indicates that most individuals fall under low to moderate income brackets. Outliers such as individuals with very high income need to be accounted during model training using scaling or other outlier handling techniques

**Insight-3:** Most individuals have either no or 1-2 children, also shows that individuals with large individuals are rare in this dataset

**Distribution of Categorical variables**

```
In [12]: categorical_cols = ['Marital Status', 'Education Level', 'Smoking Status', 'Physical Activity Level',
                           'Employment Status', 'Alcohol Consumption', 'Dietary Habits', 'Sleep Patterns',
                           'History of Substance Abuse', 'Family History of Depression', 'Chronic Medical Conditions']
```

```
# Iterate for categorical variables
for col in categorical_cols:
    plt.figure(figsize=(5, 4))
    sns.countplot(data=df, x=col)
    plt.title(f'Distribution of {col}')
    plt.xticks(rotation=45)
    plt.show()
```



**Insight-4:** Marital status is an important factor for predicting mental illness and this dataset is dominated by individuals who are married

**Insight-5:** Education may play a vital role in predicting mental, and there is a significant portion with bachelors degree and lower holding advanced degrees such as PhDs

**Insight-6:** Majority of the individuals in this dataset are non-smokers and often smoking status is associated with mental illness

**Insight-7:** Majority of individuals in this dataset are employed and employment feature is a strong predictor of mental health and unemployment is often linked to stress and anxiety

**Checking correlation of Numerical Features**

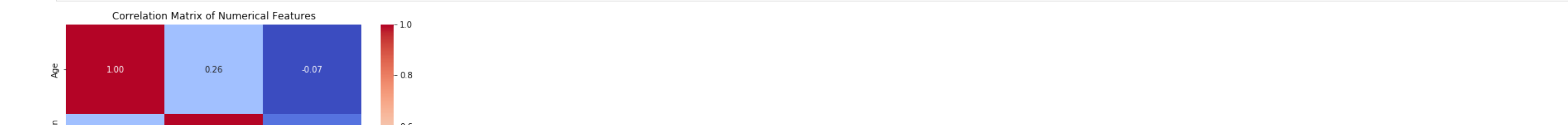
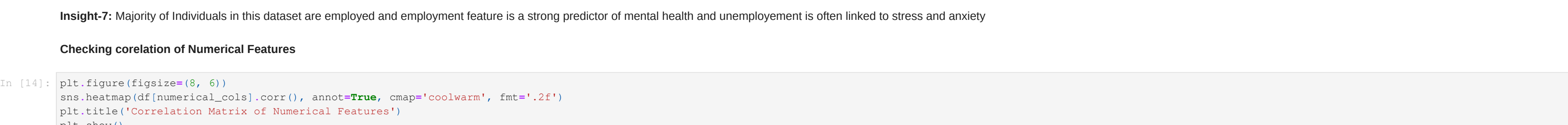
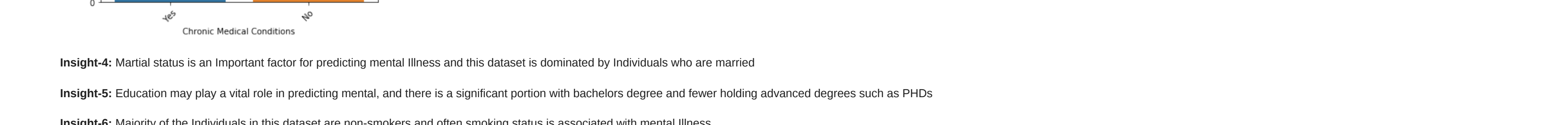
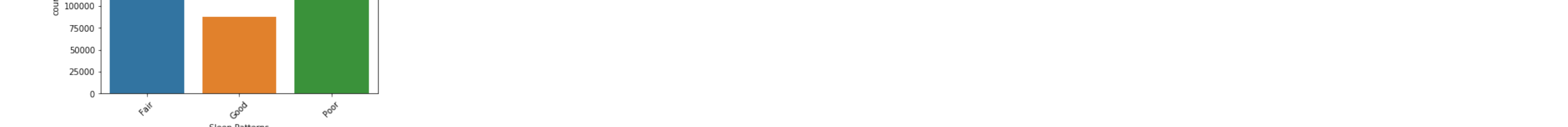
```
In [14]: plt.figure(figsize=(5, 4))
sns.heatmap(df[numerical_cols].corr(), annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Matrix of Numerical Features')
plt.show()
```



The above heatmap shows minimal correlation between numerical and categorical variables and there is no high correlation, multicollinearity should not be an issue

**Analysing relationships between categorical variables and target variables**

```
In [15]: for col in categorical_cols:
    plt.figure(figsize=(5, 4))
    sns.countplot(data=df, x=col, hue='History of Mental Illness')
    plt.title(f'{col} vs. History of Mental Illness')
    plt.xticks(rotation=45)
    plt.show()
```



**Insight-8:** This dataset portrays that married individuals seem to have lower rates of Mental health issues while single and widowed individuals show higher proportions of mental illness

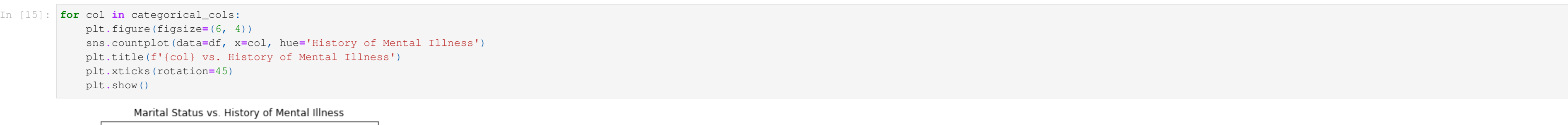
**Insight-9:** Mental illness is higher among those with lower educational backgrounds while it is lower with advanced degrees

**Insight-10:** Current and former smokers have higher rates of mental illness compared to non-smokers.

**Insight-11:** Employment may be one of the strongest predictors as no employment often leads to major stress problems

**Checking for Outliers in the Numerical variables**

```
In [19]: for col in numerical_cols:
    plt.figure(figsize=(5, 4))
    sns.boxplot(data=df, x=col)
    plt.title(f'Boxplot of {col}')
    plt.show()
```



**Age:** No significant outliers found with most of them falling between 20 and 80

**Income:** Small number of individuals earning more than the majority. These outliers may need to be removed

**Number of Children:** No significant outliers found as the maximum is capped at 4

**Running Chi-Square Test to understand Feature Importance**

```
In [20]: from sklearn.feature_selection import chi2
from sklearn.preprocessing import StandardScaler
import pandas as pd

# List of categorical columns
categorical_cols = ['Marital Status', 'Education Level', 'Smoking Status', 'Physical Activity Level',
                   'Employment Status', 'Alcohol Consumption', 'Dietary Habits', 'Sleep Patterns',
                   'History of Substance Abuse', 'Family History of Depression', 'Chronic Medical Conditions']

# Encoding the target variable
df['History of Mental Illness'] = df['History of Mental Illness'].map({'Yes': 1, 'No': 0})

# Encoding the categorical columns
df_encoded = pd.get_dummies(df, columns=categorical_cols, drop_first=False)

# Features (X) and target (y)
X = df_encoded.drop(columns='History of Mental Illness')
y = df_encoded['History of Mental Illness']

# Dropping irrelevant 'Name' column
X = X.drop('Name', axis=1)

# Normalizing data for Chi-Square test (Chi-Square test requires non-negative values)
X_scaled = X.copy().astype(float)
X_scaled = X_scaled * 100000

# Performing chi-square test
chi2_scores = chi2(X_scaled, y)

# Display p-values and Feature Importance
p_values = pd.Series(chi2_scores, index=X_scaled.columns)
significant_features = p_values[p_values < 0.05] # Features significant at p < 0.05
print(f'Important features based on Chi-square test: {significant_features}')

# Important features based on Chi-square test
Important_features = [
    'Marital Status_Divorced',
    'Marital Status_Single',
    'Marital Status_Widowed',
    'Education Level_Associate Degree',
    'Education Level_Bachelor's Degree',
    'Education Level_High School',
    'Education Level_Master's Degree',
    'Education Level_PhD',
    'Smoking Status_Current',
    'Physical Activity Level_Sedentary',
    'Physical Activity Level_Active',
    'Employment Status_Unemployed',
    'Alcohol Consumption_High',
    'Dietary Habits_Unhealthy',
    'Sleep Patterns_Fair',
    'Family History of Depression_Yes',
    'Chronic Medical Conditions_No',
    'Chronic Medical Conditions_Yes'
]
```

Key Insights

- The following Chi-Square Test reveals that Socio-demographic variables (age, marital status, education), lifestyle choices (physical activity, smoking and alcohol consumption) and health related factors (sleep patterns, family history of depression, chronic medical conditions) are associated with mental health

In [21]: