

21BDS0379

SHREYASHA SHRESTHA

EDA THEORY- D1

Digital Assignment 1

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# Load dataset
df = pd.read_csv('MyDataset_EDA.csv')
```

```
# Display basic info
df.info()
print("\nSummary Statistics:")
print(df.describe(include='all'))
```

```
32 G3 395 non-null int64
dtypes: int64(16), object(17)
memory usage: 102.0+ KB
```

Summary Statistics:

	school	sex	age	address	famsize	Pstatus	Medu	\
count	395	393	395.000000	394	394	395	395.000000	
unique	2	2	NaN	3	2	3	NaN	
top	GP	F	NaN	U	GT3	T	NaN	
freq	349	208	NaN	304	280	348	NaN	
mean	NaN	NaN	16.696203	NaN	NaN	NaN	2.749367	
std	NaN	NaN	1.276043	NaN	NaN	NaN	1.094735	
min	NaN	NaN	15.000000	NaN	NaN	NaN	0.000000	
25%	NaN	NaN	16.000000	NaN	NaN	NaN	2.000000	
50%	NaN	NaN	17.000000	NaN	NaN	NaN	3.000000	
75%	NaN	NaN	18.000000	NaN	NaN	NaN	4.000000	
max	NaN	NaN	22.000000	NaN	NaN	NaN	4.000000	

	Fedu	Mjob	Fjob	...	famrel	freetime	goout	\
count	395.000000	395	395	...	395.000000	395.000000	395.000000	
unique	NaN	5	5	...	NaN	NaN	NaN	
top	NaN	other	other	...	NaN	NaN	NaN	
freq	NaN	141	217	...	NaN	NaN	NaN	
mean	2.521519	NaN	NaN	...	3.944304	3.235443	3.108861	
std	1.088201	NaN	NaN	...	0.896659	0.998862	1.113278	
min	0.000000	NaN	NaN	...	1.000000	1.000000	1.000000	
25%	2.000000	NaN	NaN	...	4.000000	3.000000	2.000000	
50%	2.000000	NaN	NaN	...	4.000000	3.000000	3.000000	
75%	3.000000	NaN	NaN	...	5.000000	4.000000	4.000000	
max	4.000000	NaN	NaN	...	5.000000	5.000000	5.000000	

	Dalc	Walc	health	absences	G1	\
count	395.000000	395.000000	395.000000	395.000000	395.000000	
unique	NaN	NaN	NaN	NaN	NaN	
top	NaN	NaN	NaN	NaN	NaN	
freq	NaN	NaN	NaN	NaN	NaN	
mean	1.481013	2.291139	3.554430	5.708861	10.908861	
std	0.890741	1.287897	1.390303	8.003096	3.319195	
min	1.000000	1.000000	1.000000	0.000000	3.000000	
25%	1.000000	1.000000	3.000000	0.000000	8.000000	
50%	1.000000	2.000000	4.000000	4.000000	11.000000	
75%	2.000000	3.000000	5.000000	8.000000	13.000000	
max	5.000000	5.000000	5.000000	75.000000	19.000000	

	G2	G3
count	395.000000	395.000000
unique	NaN	NaN
top	NaN	NaN
freq	NaN	NaN
mean	10.713924	10.415190
std	3.761505	4.581443
min	0.000000	0.000000
25%	9.000000	8.000000
50%	11.000000	11.000000
75%	13.000000	14.000000
max	19.000000	20.000000

[11 rows x 33 columns]


```
# Handling outliers (Removing values beyond 1.5*IQR)
def remove_outliers(df):
    for col in ['absences']:
        Q1 = df[col].quantile(0.25)
        Q3 = df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
        df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]
    return df

df = remove_outliers(df)

# Encoding categorical variables
for col in df.select_dtypes(include=['object']).columns:
    df[col] = df[col].astype('category').cat.codes
```

Univariate Analysis

```
# Central Tendency Measures
for col in df.select_dtypes(include=['float64', 'int64']).columns:
    print(f"\n{col}:")
    print(f"Mean: {df[col].mean()}")
    print(f"Median: {df[col].median()}")
    print(f"Mode: {df[col].mode()[0]}")


# Dispersion Measures
for col in df.select_dtypes(include=['float64', 'int64']).columns:
    print(f"\n{col}:")
    print(f"Range: {df[col].max() - df[col].min()}")
    print(f"Variance: {df[col].var()}")
    print(f"Standard Deviation: {df[col].std()}")
    print(f"Min: {df[col].min()}")
    print(f"Max: {df[col].max()}")
    print(f"Q1: {df[col].quantile(0.25)}")
    print(f"Q3: {df[col].quantile(0.75)}")
    print(f"IQR: {df[col].quantile(0.75) - df[col].quantile(0.25)}")
```



```
Min: 0  
Max: 19  
Q1: 9.0  
Q3: 13.0  
IQR: 4.0
```

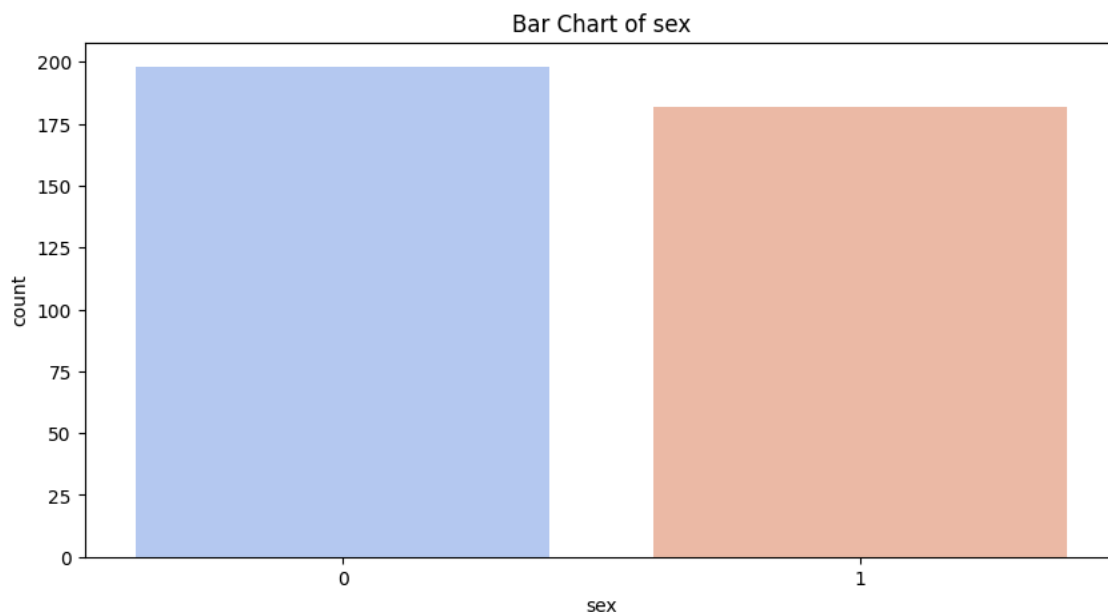
```
G3:  
Range: 20  
Variance: 21.27343424524371  
Standard Deviation: 4.612313329040397  
Min: 0  
Max: 20  
Q1: 8.0  
Q3: 14.0  
IQR: 6.0
```

```
categorical_cols = ['sex', 'address', 'famsize', 'schoolsup', 'age_group']  
for cat_col in categorical_cols:  
    if cat_col in df.columns:  
        plt.figure(figsize=(10,5))  
        sns.countplot(x=df[cat_col].astype(str), palette='coolwarm')  
        plt.title(f'Bar Chart of {cat_col}')  
        plt.show()  
    else:  
        print(f"Column '{cat_col}' not found in DataFrame, skipping.")
```

 <ipython-input-37-c8f26cb134f5>:5: FutureWarning:

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

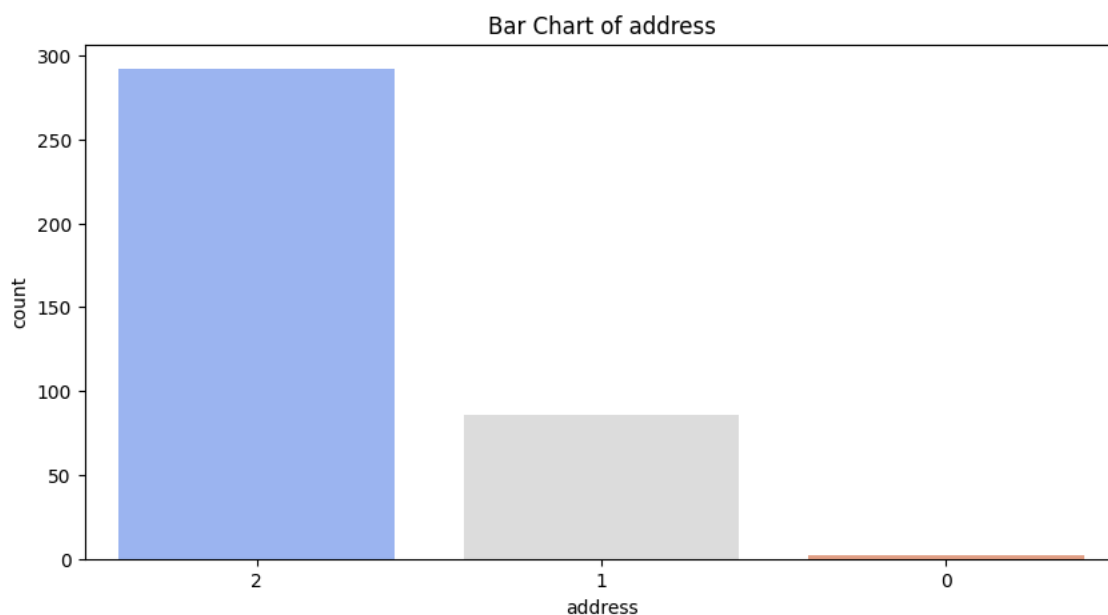
```
sns.countplot(x=df[cat_col].astype(str), palette='coolwarm')
```



<ipython-input-37-c8f26cb134f5>:5: FutureWarning:

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

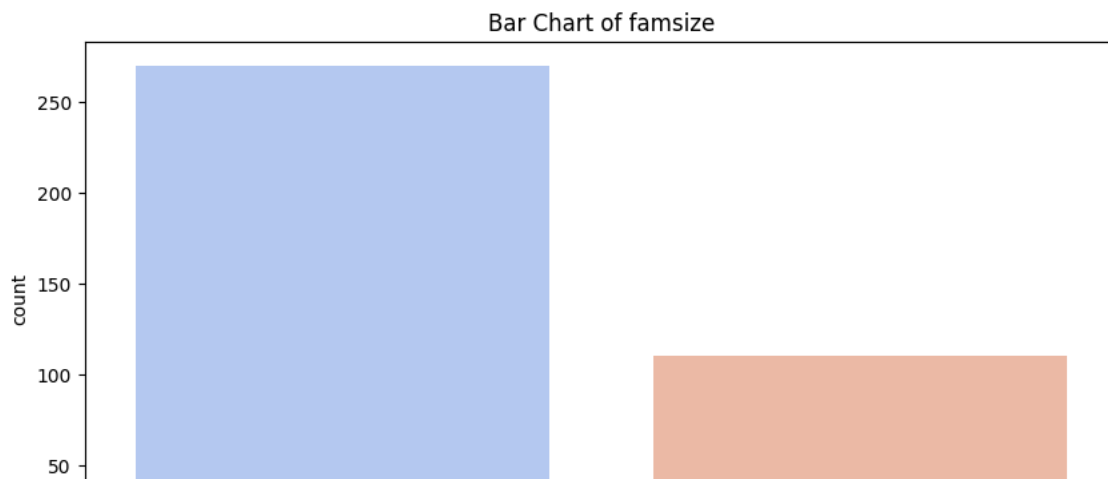
```
sns.countplot(x=df[cat_col].astype(str), palette='coolwarm')
```

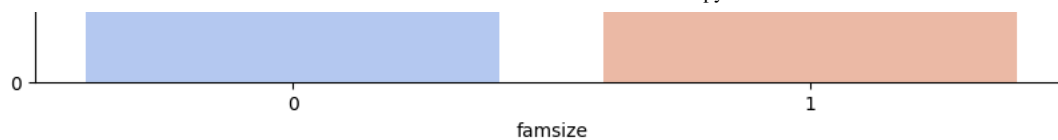


<ipython-input-37-c8f26cb134f5>:5: FutureWarning:

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

```
sns.countplot(x=df[cat_col].astype(str), palette='coolwarm')
```



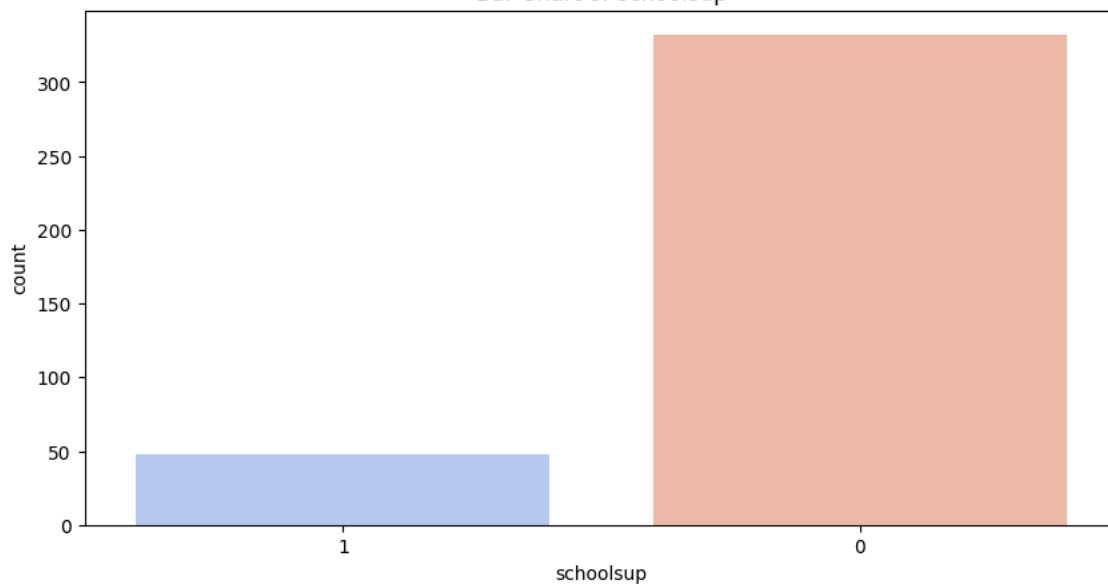


<ipython-input-37-c8f26cb134f5>:5: FutureWarning:

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

```
sns.countplot(x=df[cat_col].astype(str), palette='coolwarm')
```

Bar Chart of schoolsup

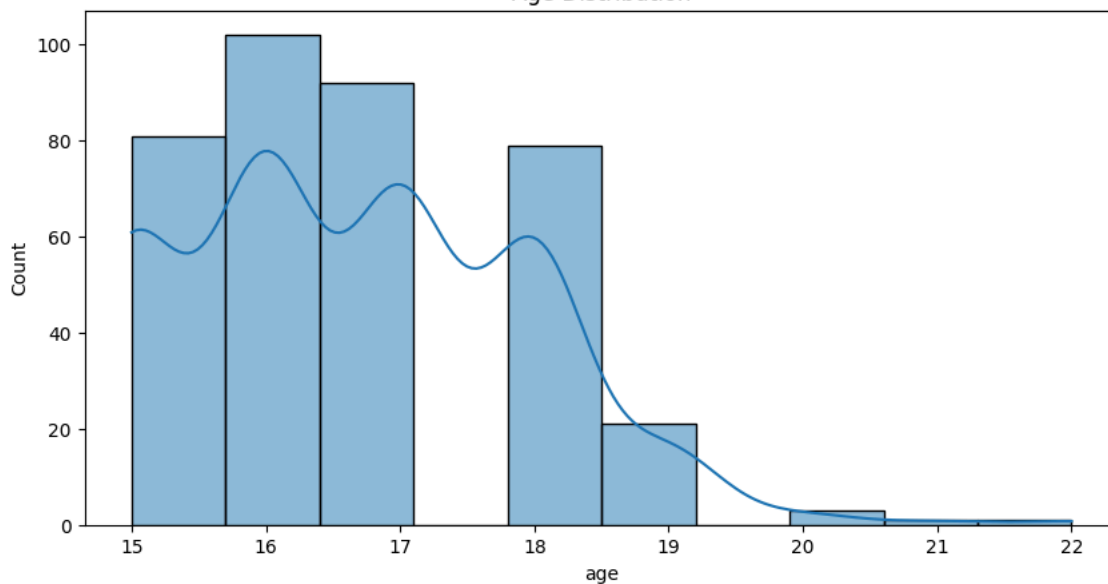


Column 'age_group' not found in DataFrame, skipping.

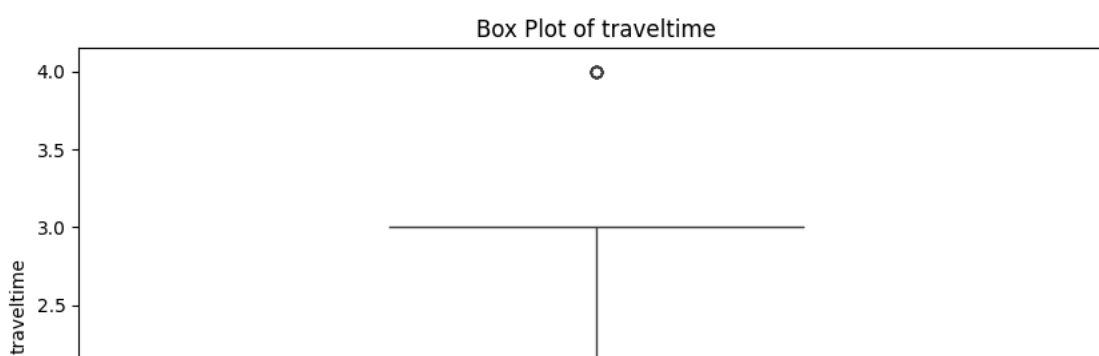
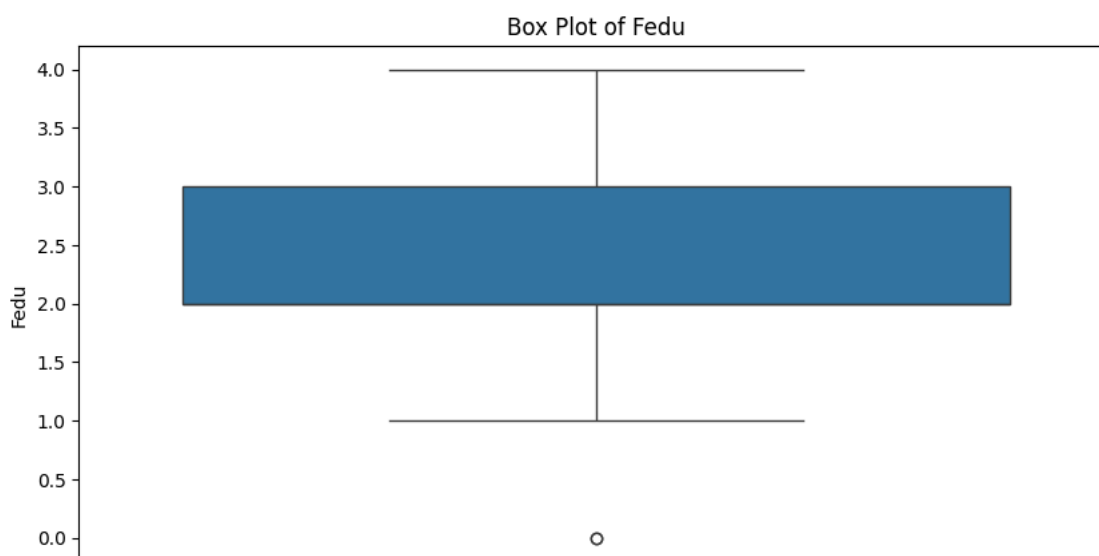
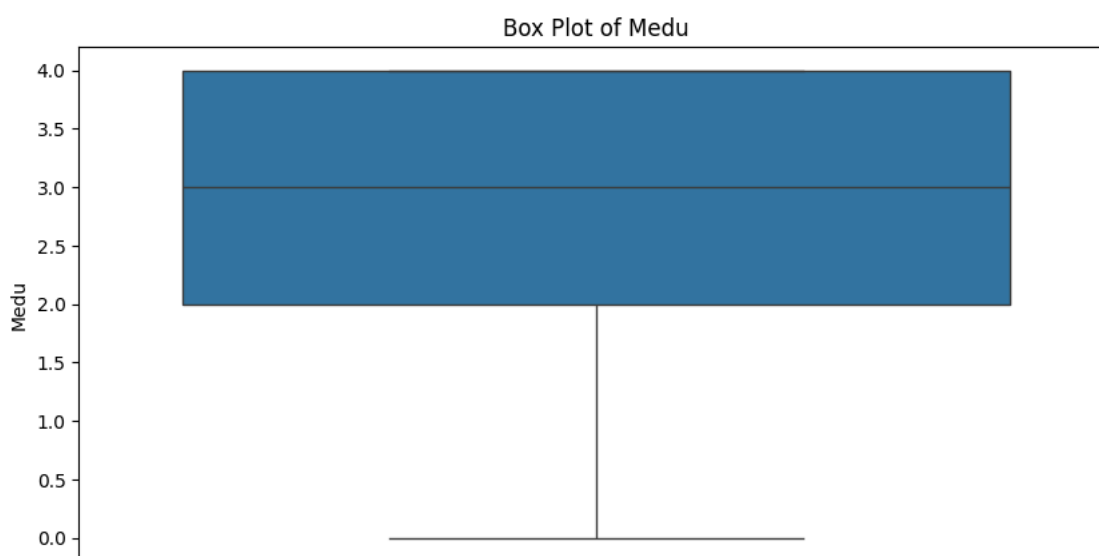
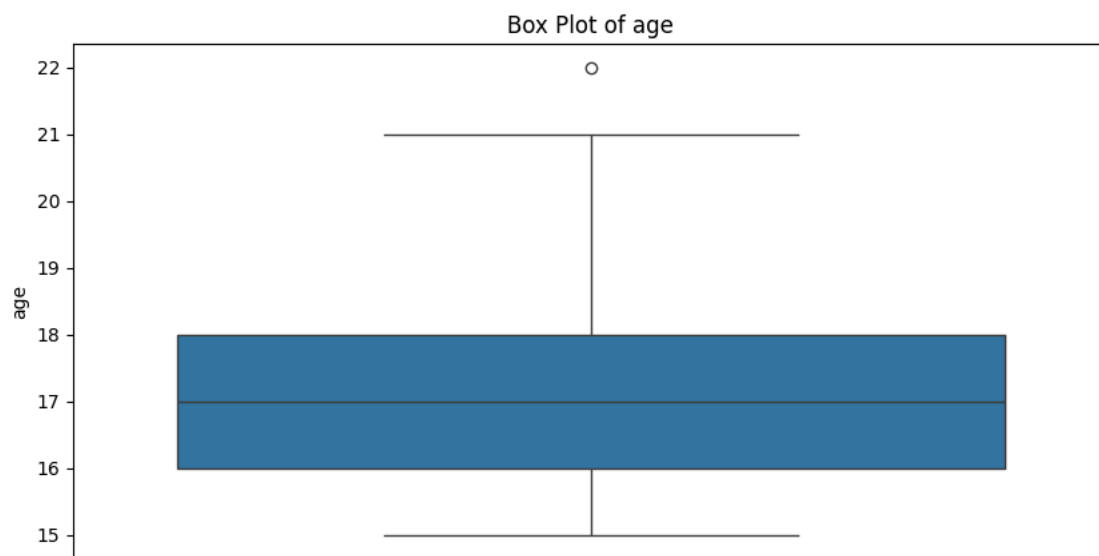
```
plt.figure(figsize=(10,5))
sns.histplot(df['age'], bins=10, kde=True)
plt.title('Age Distribution')
plt.show()
```

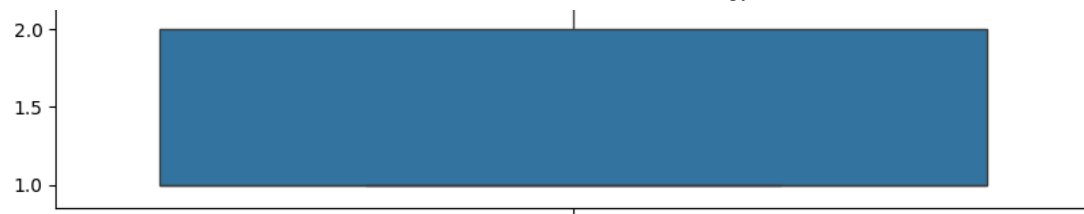


Age Distribution

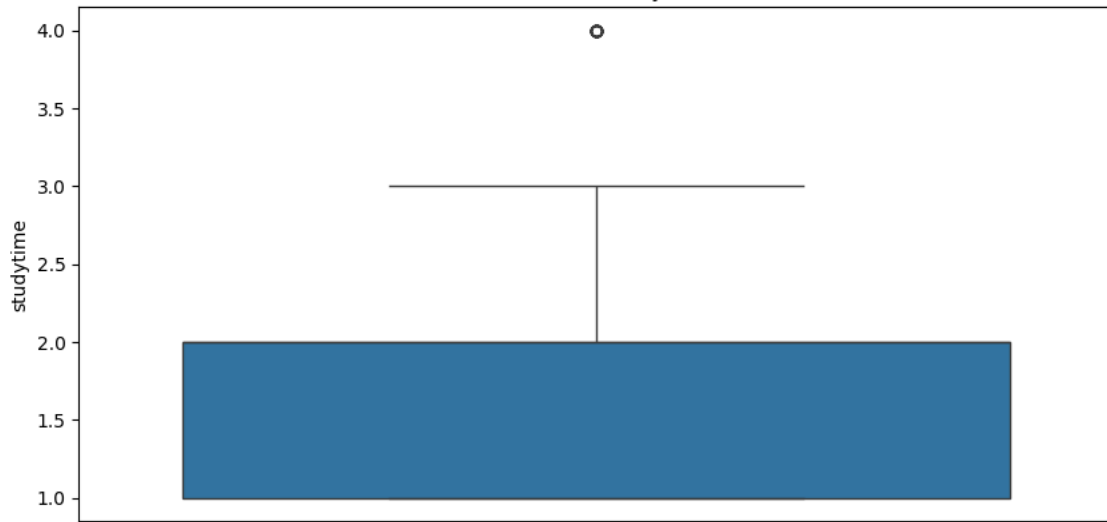


```
# Box plot for numeric columns
for col in df.select_dtypes(include=['float64', 'int64']).columns:
    plt.figure(figsize=(10,5))
    sns.boxplot(y=df[col])
    plt.title(f'Box Plot of {col}')
    plt.show()
```

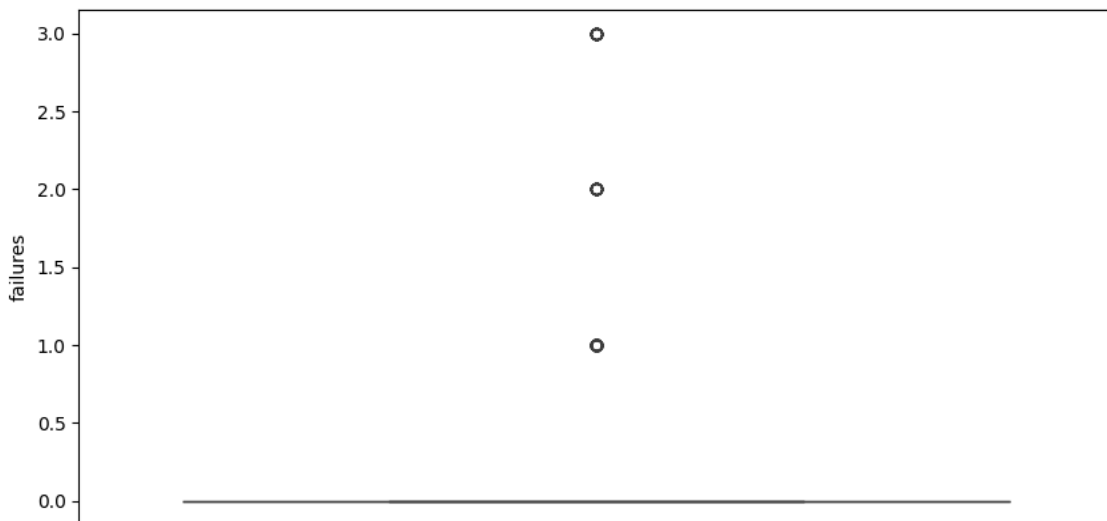




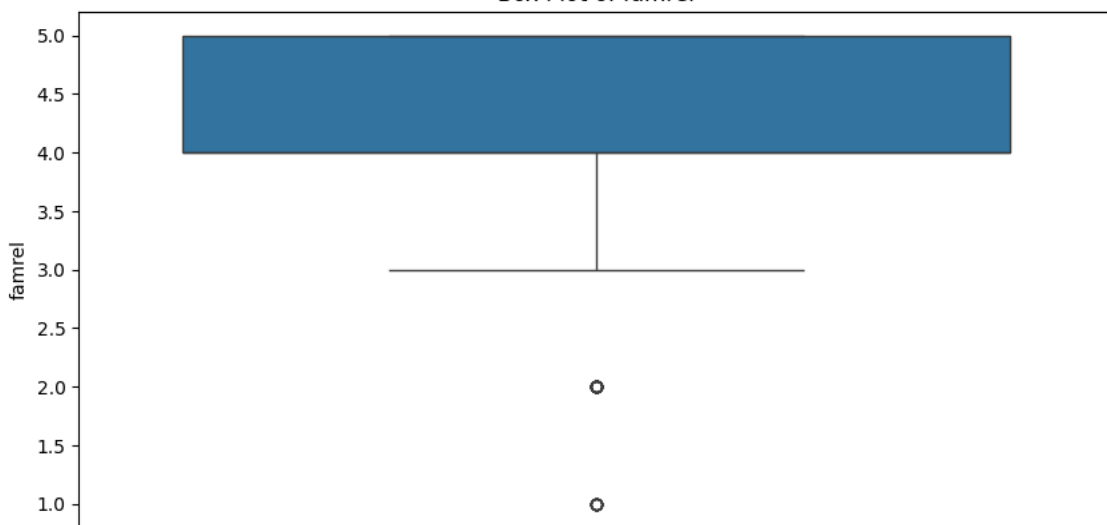
Box Plot of studytime



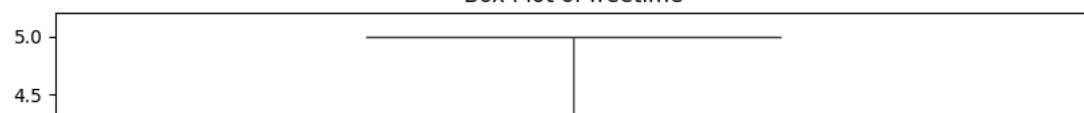
Box Plot of failures

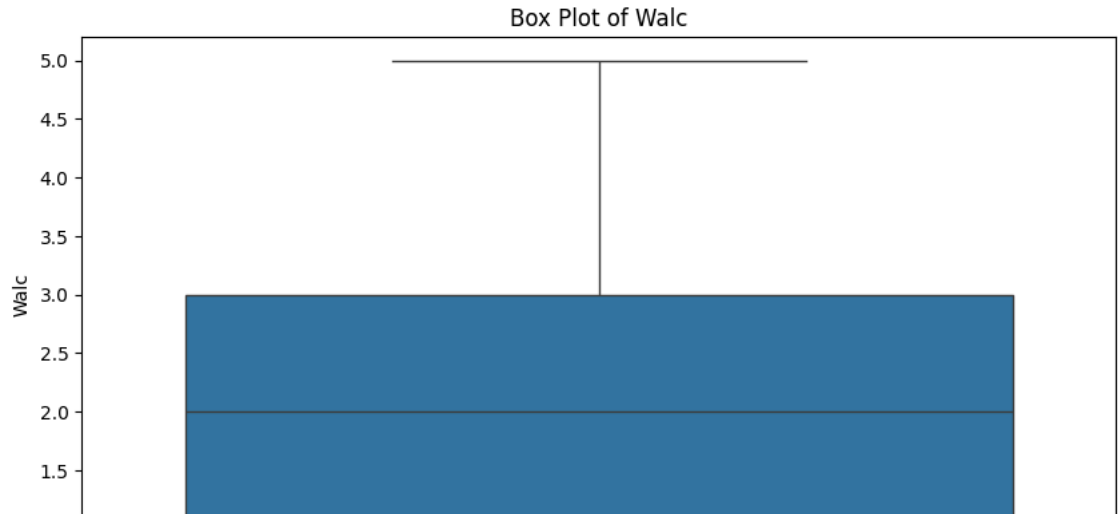
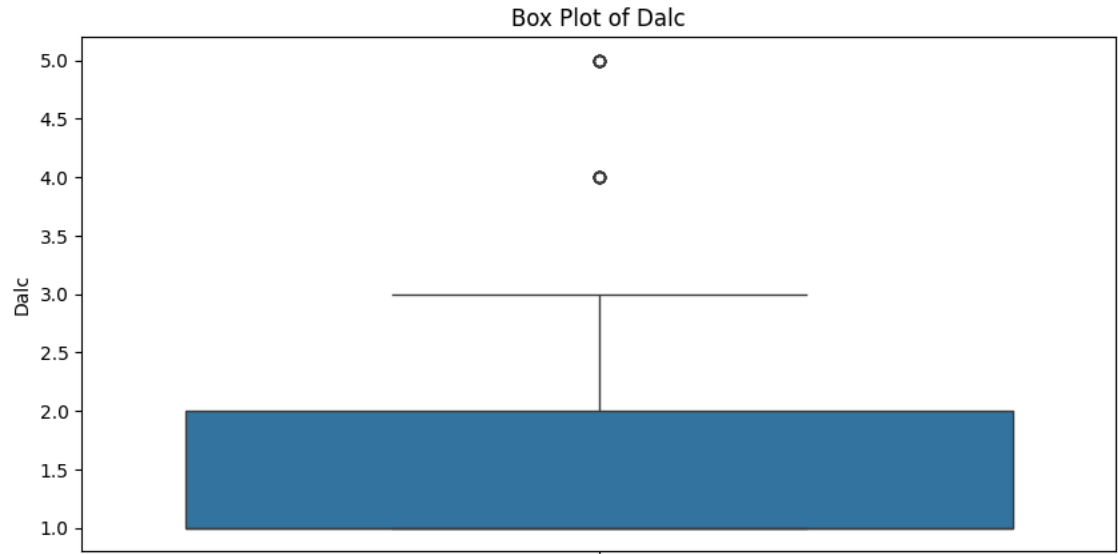
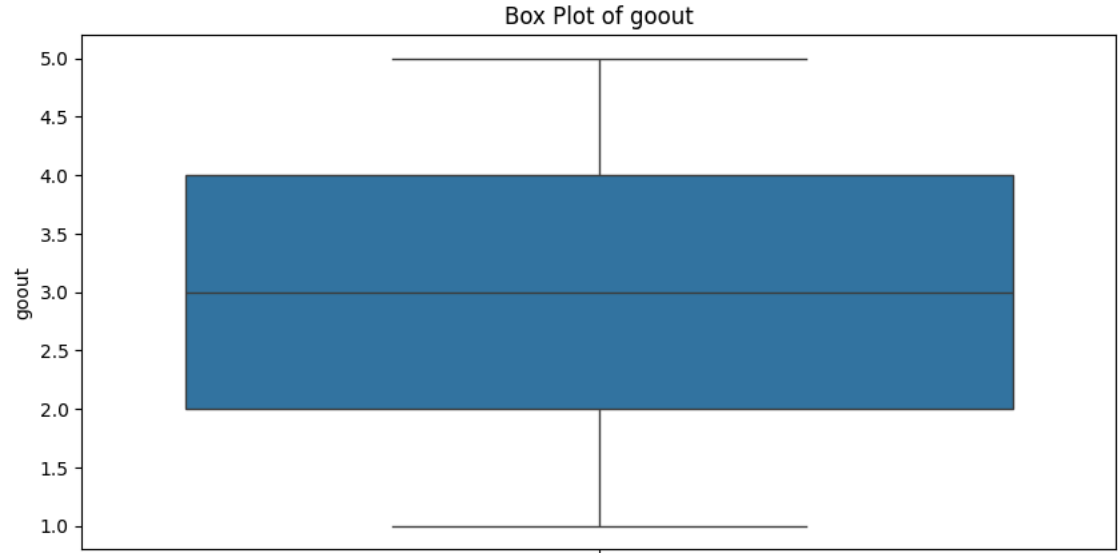
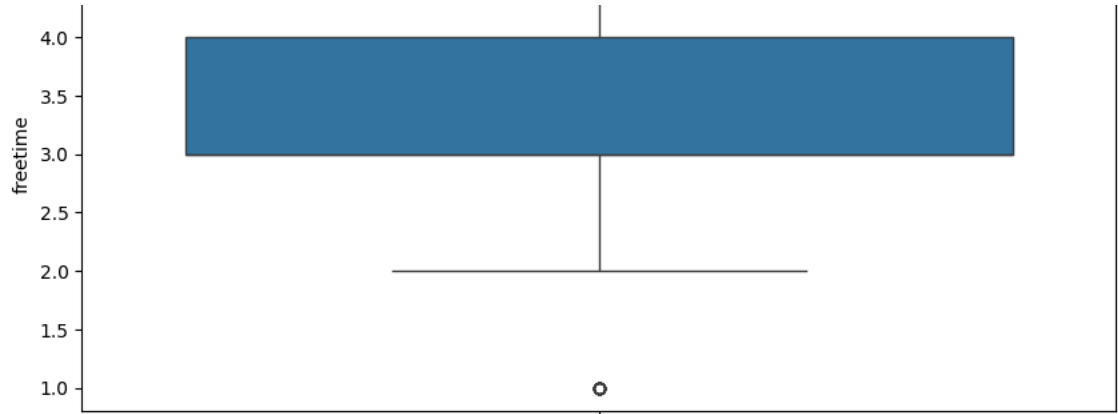


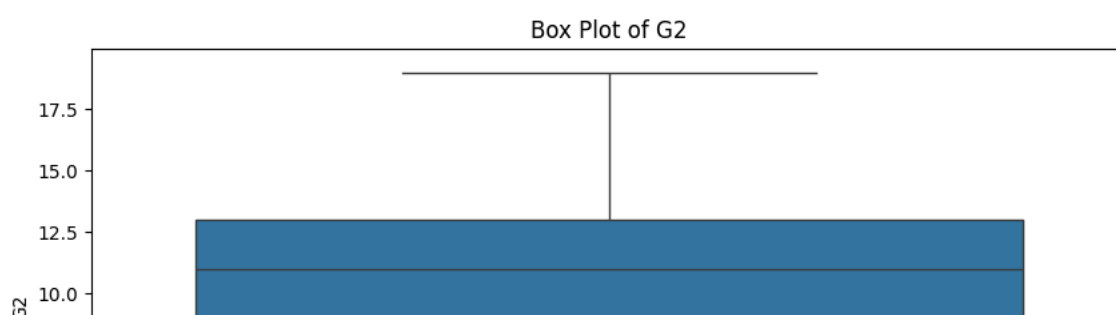
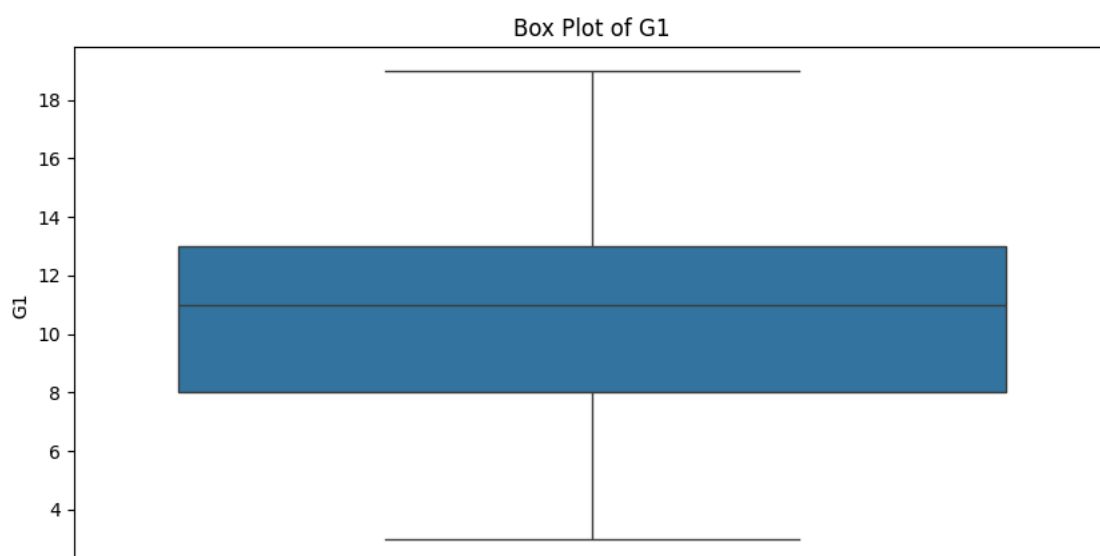
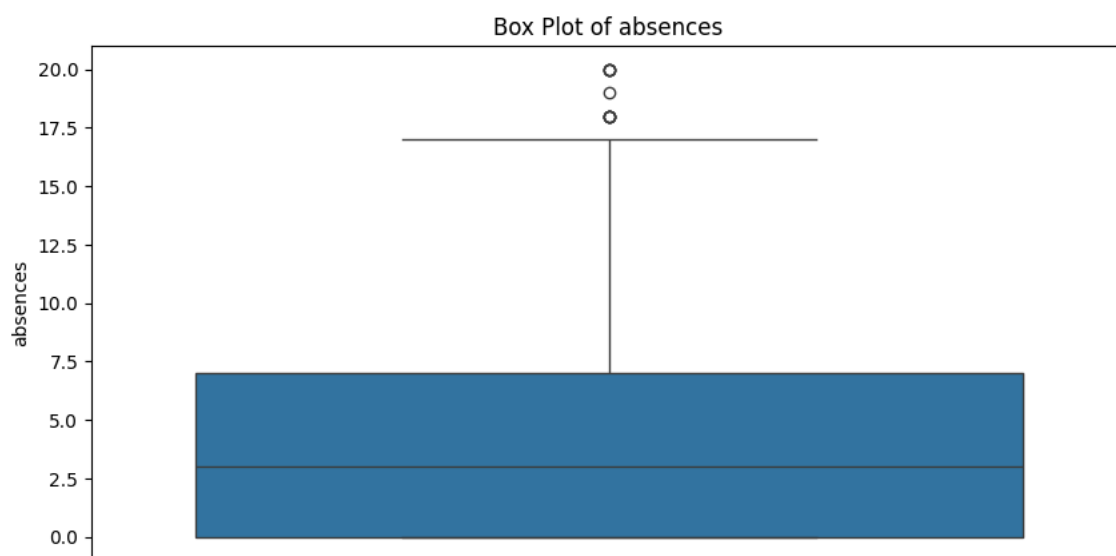
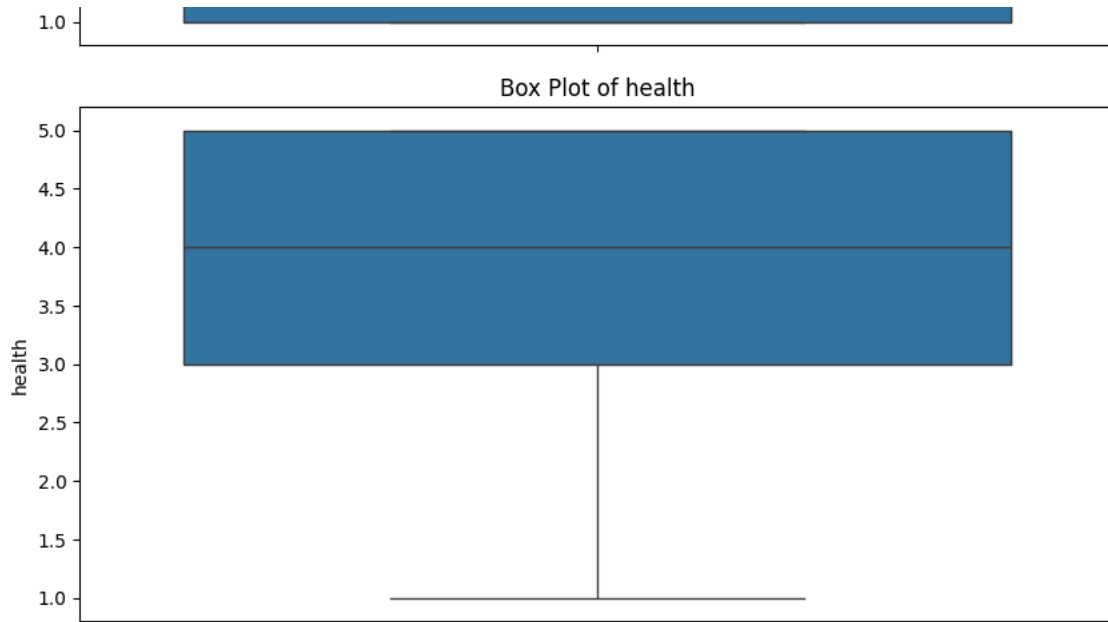
Box Plot of famrel

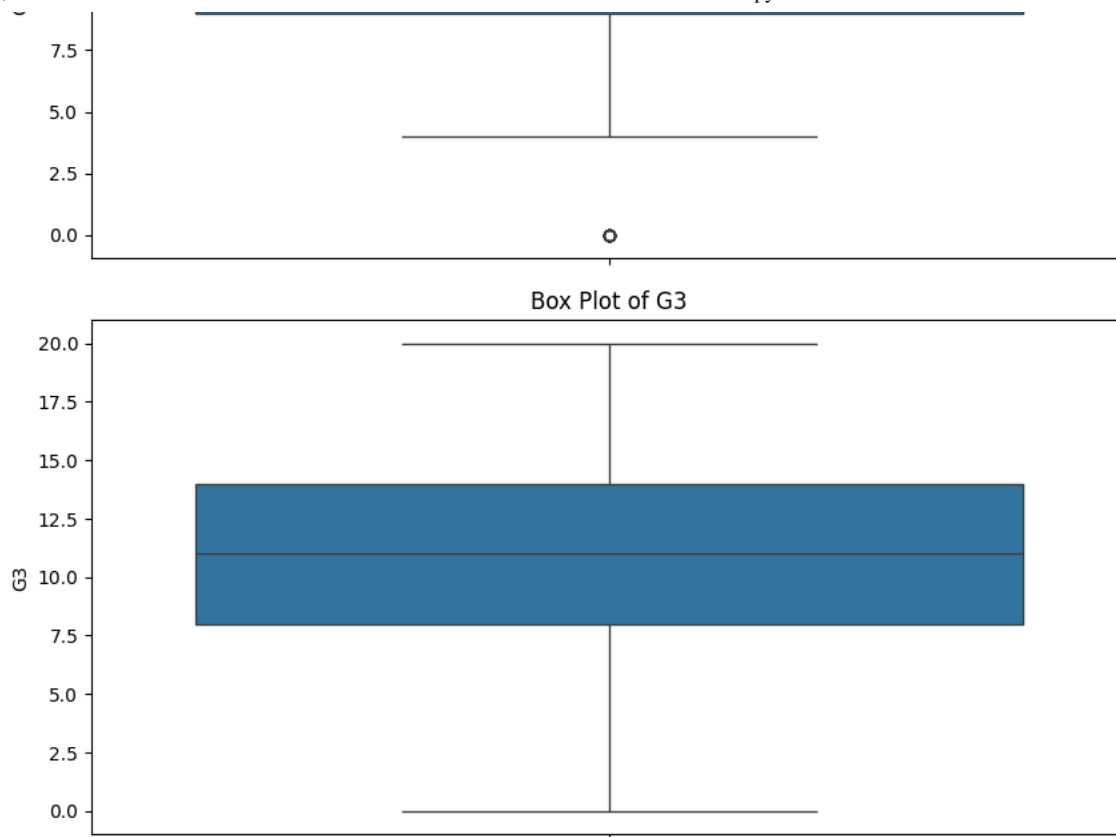


Box Plot of freetime





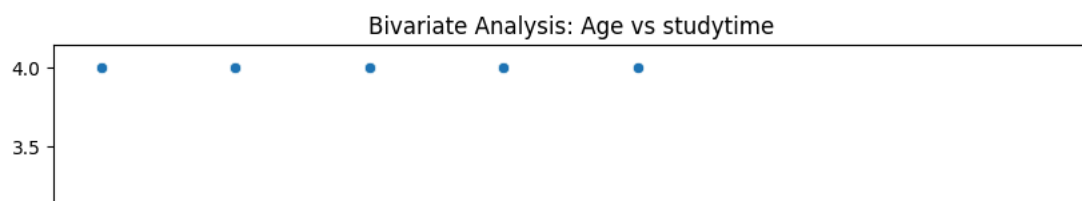
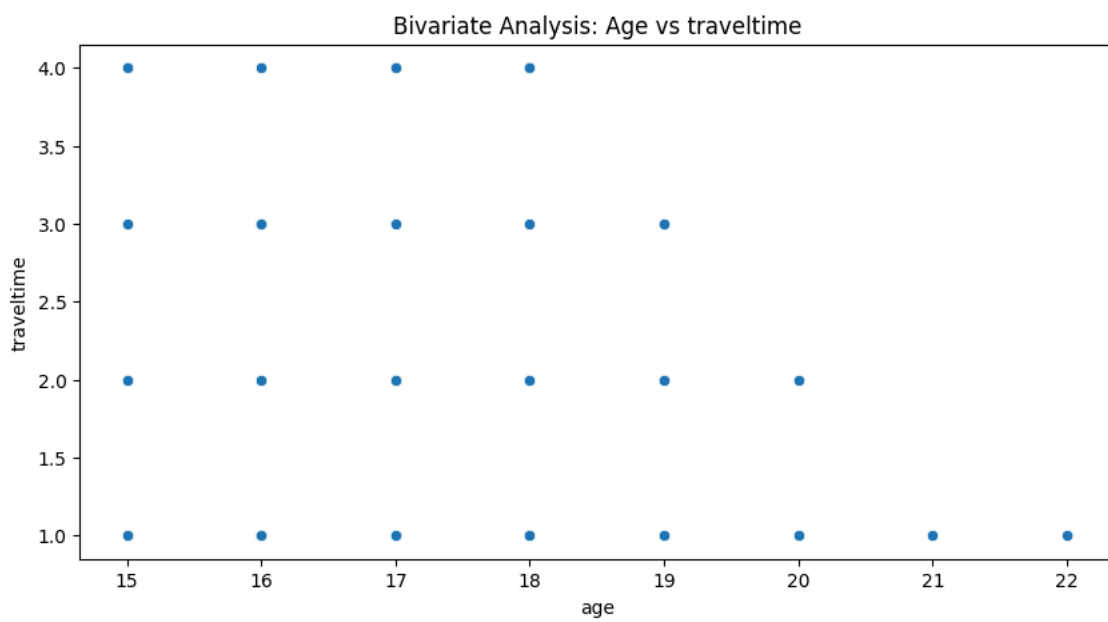
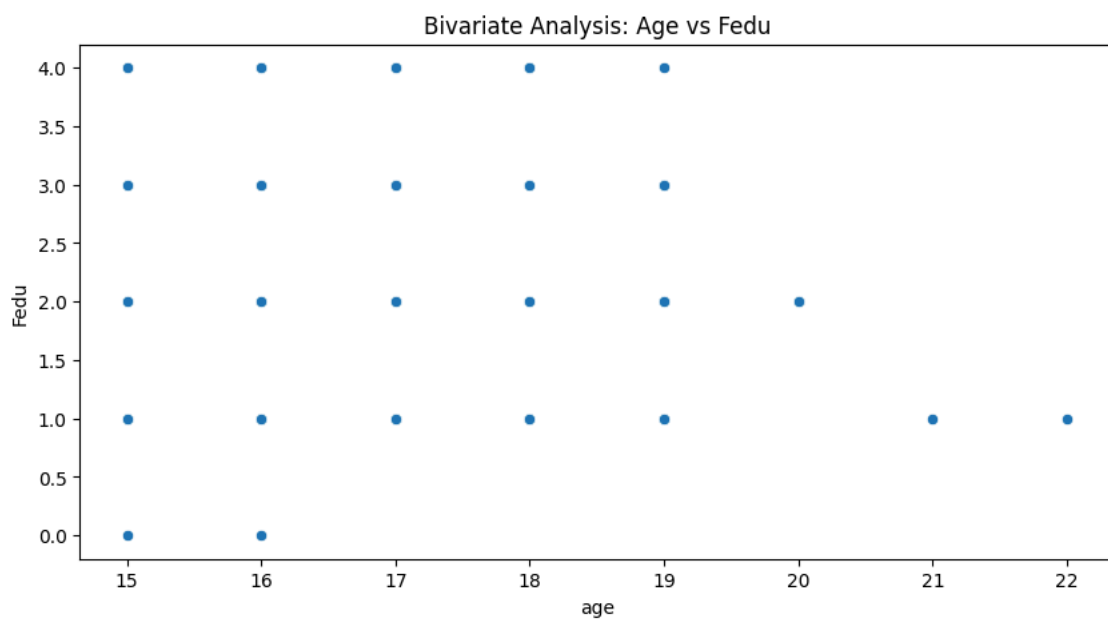
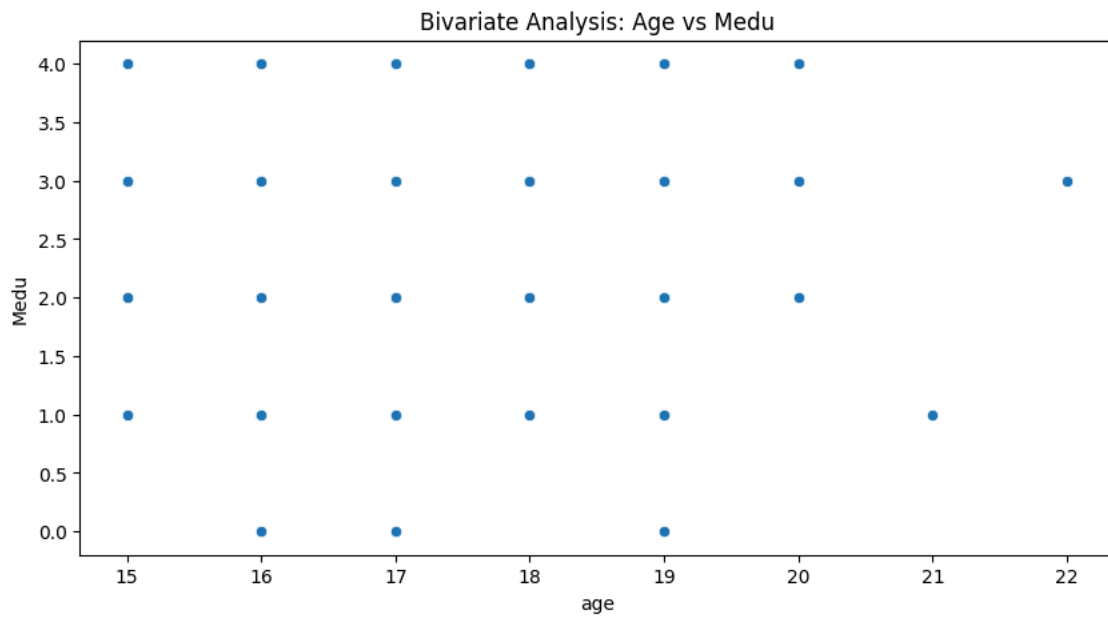


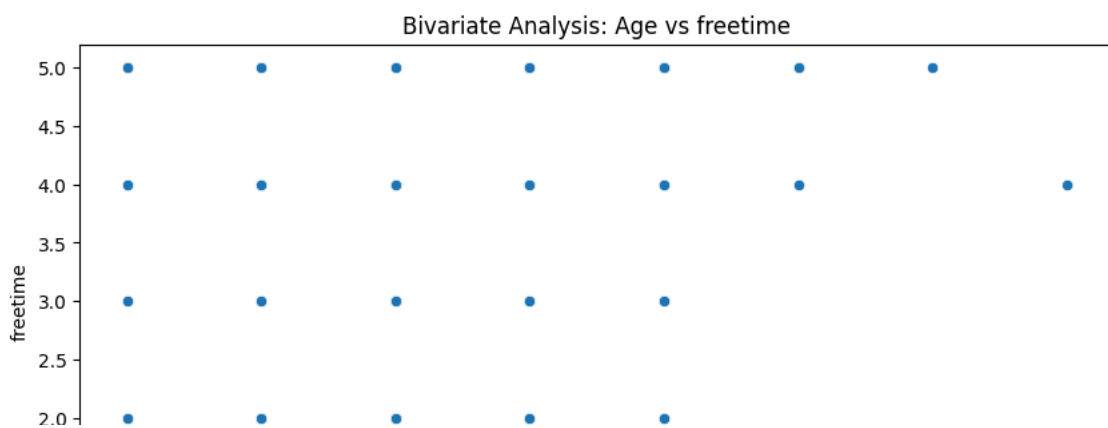
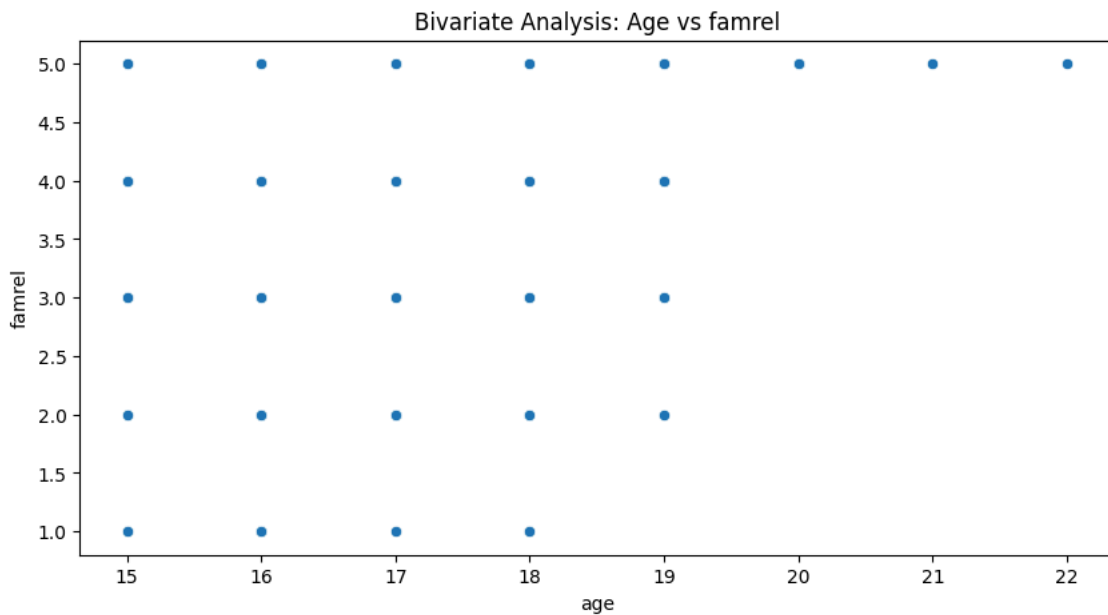
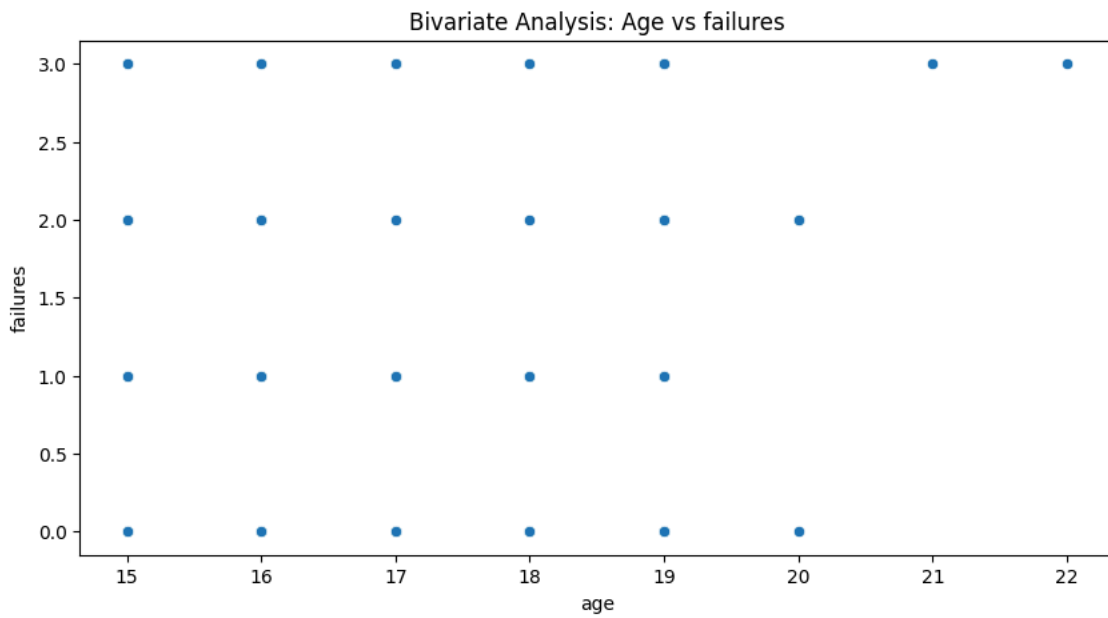
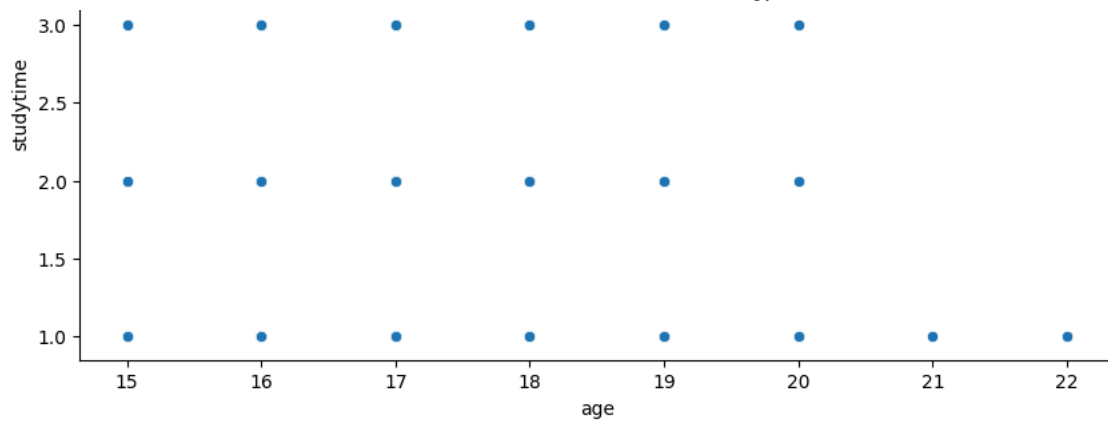


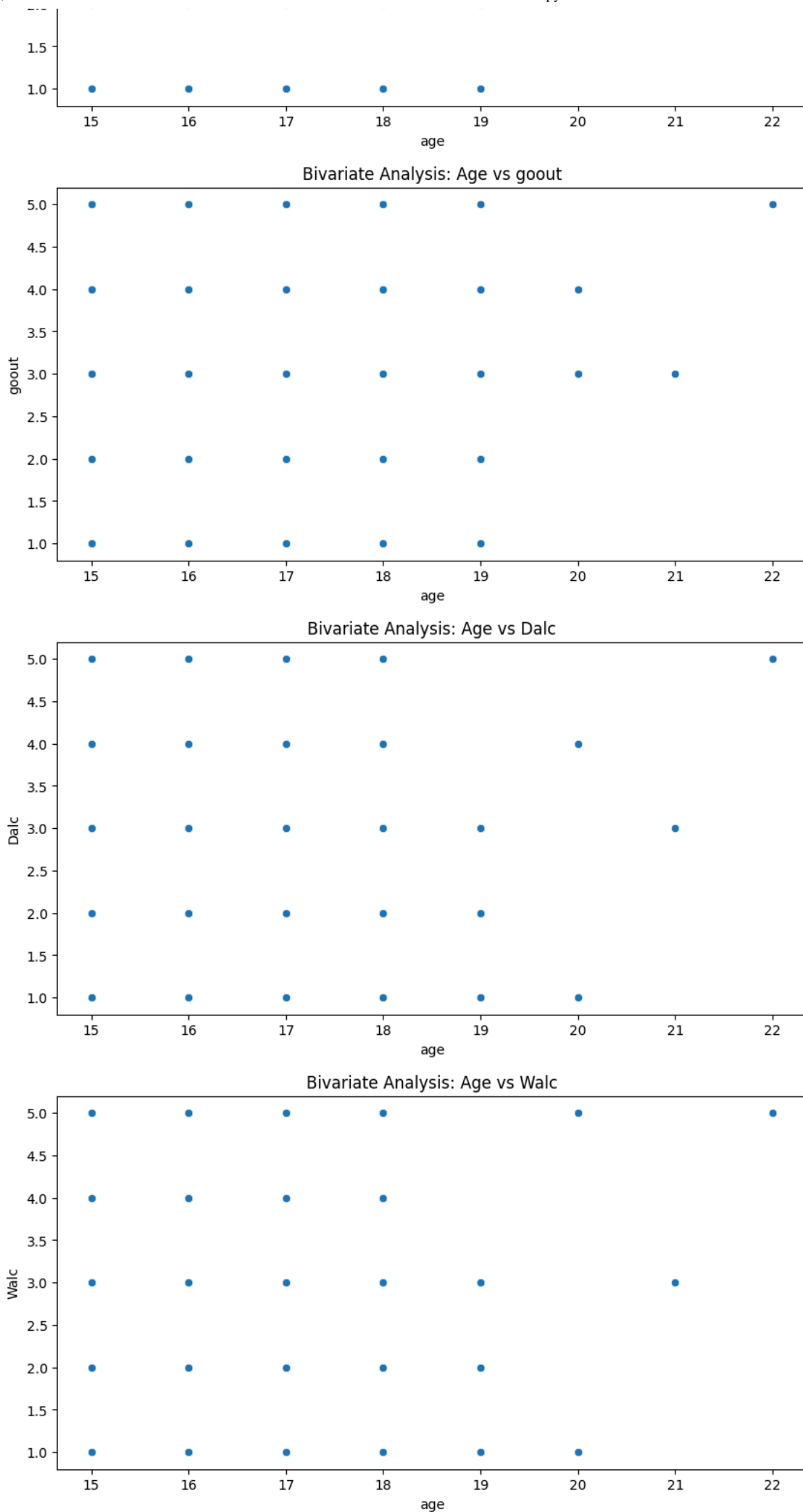
Double-click (or enter) to edit

Bivariate Analysis

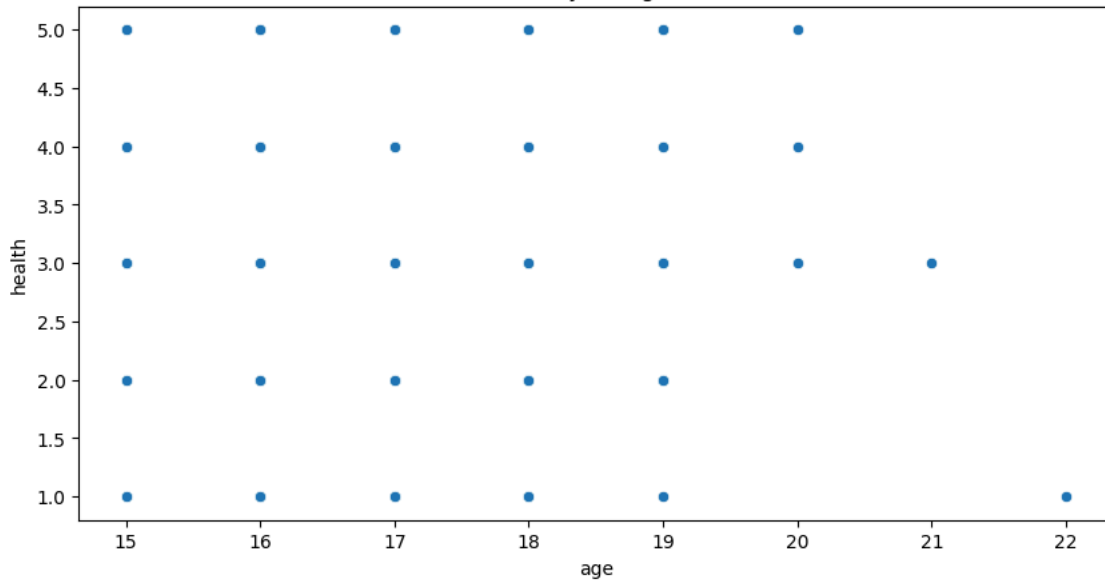
```
#Bivariate Analysis
for col in df.select_dtypes(include=['float64', 'int64']).columns:
    if col != 'age': # Avoid redundant comparison
        plt.figure(figsize=(10,5))
        sns.scatterplot(x=df['age'], y=df[col])
        plt.title(f'Bivariate Analysis: Age vs {col}')
        plt.show()
```



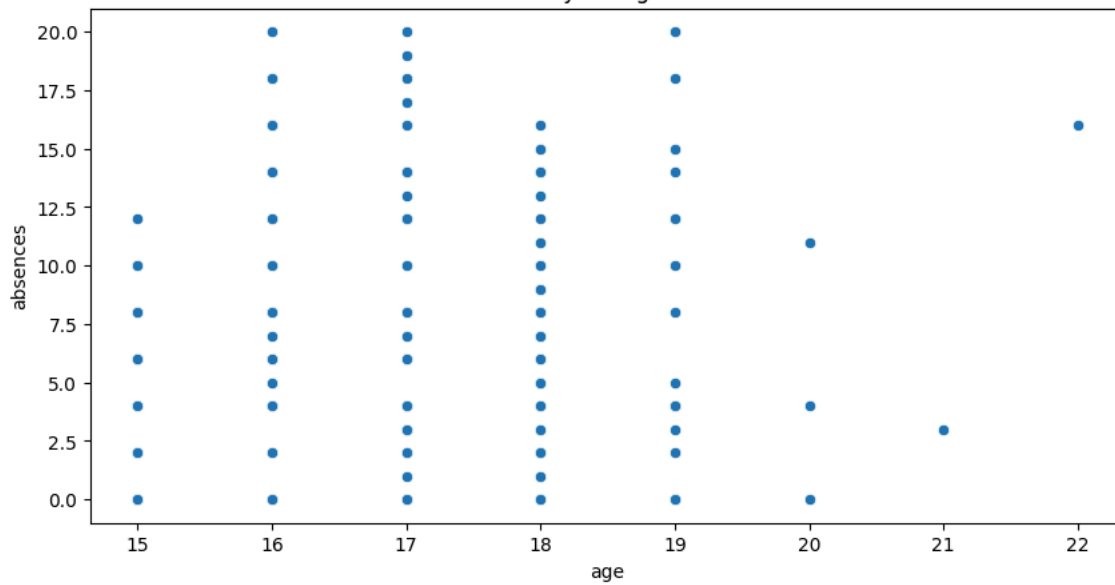




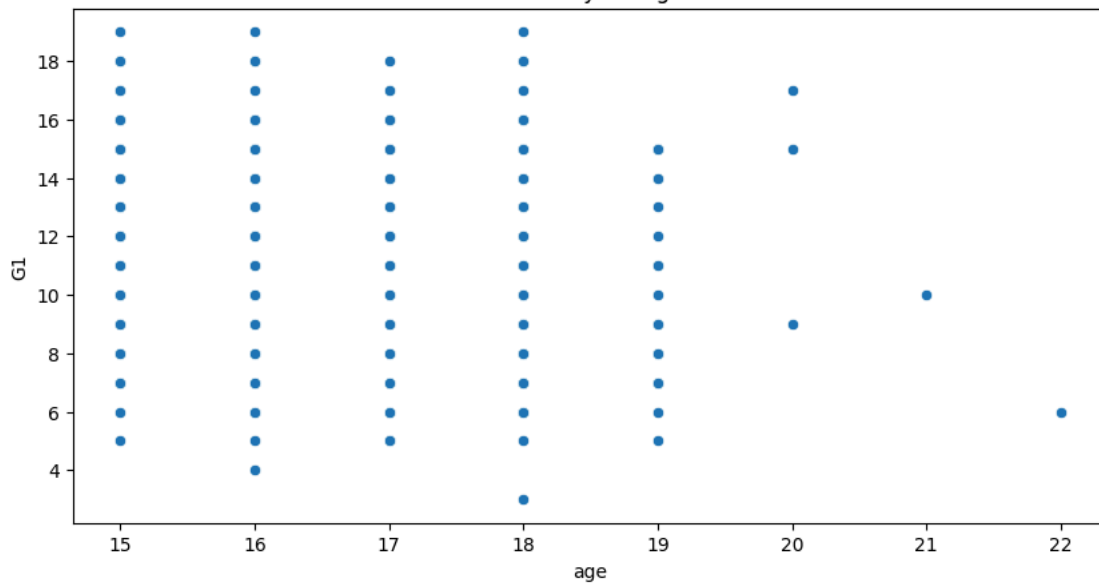
Bivariate Analysis: Age vs health



Bivariate Analysis: Age vs absences

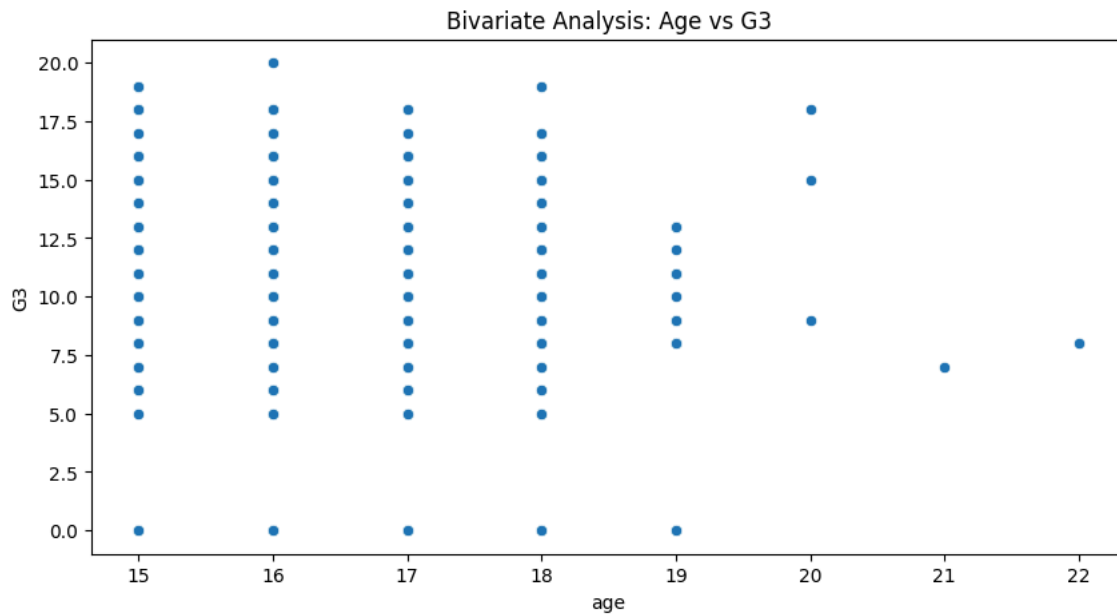
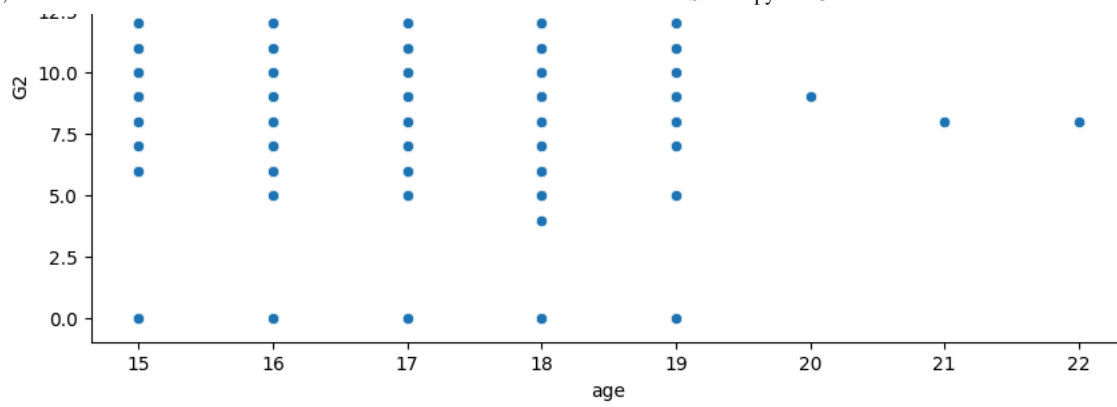


Bivariate Analysis: Age vs G1

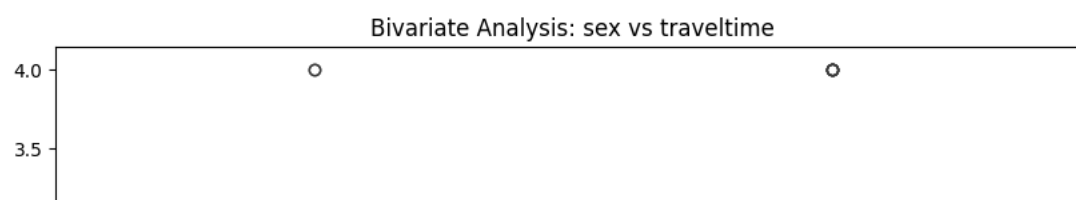
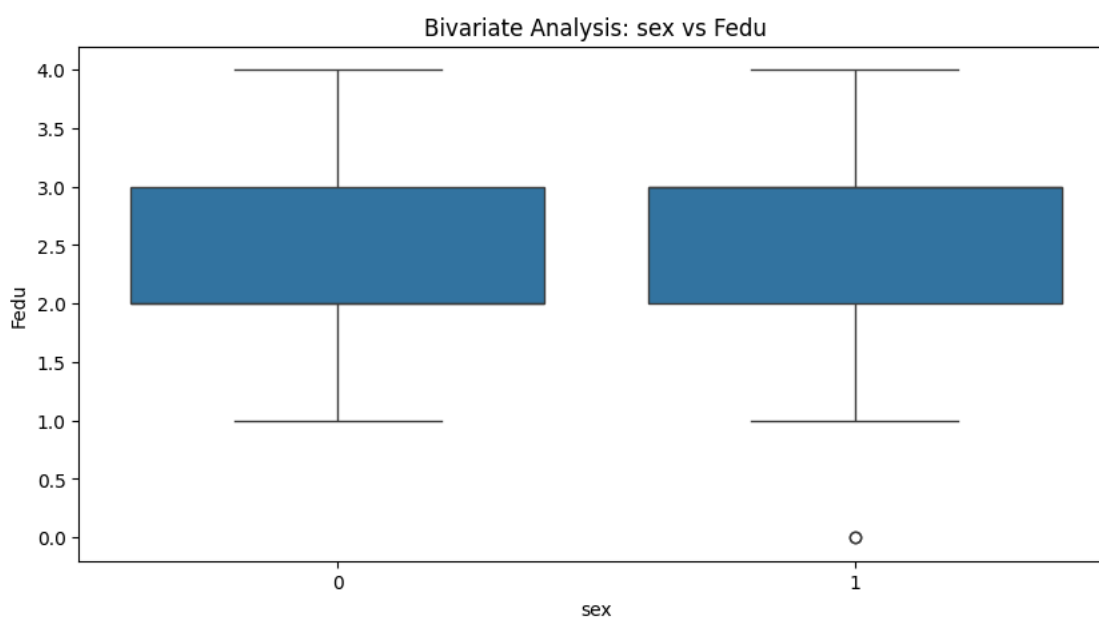
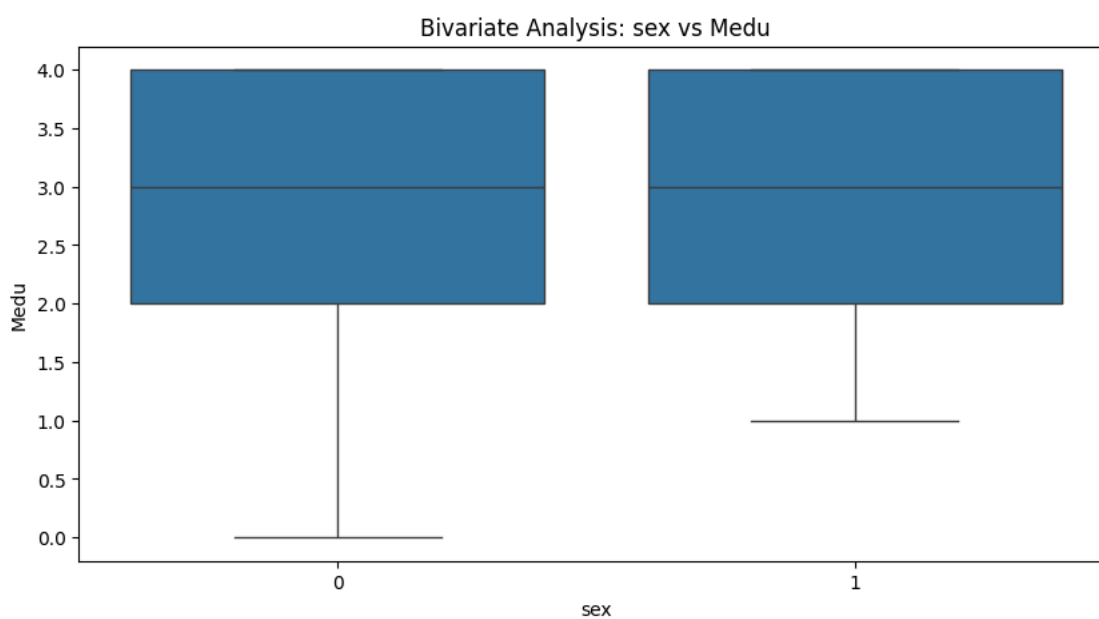
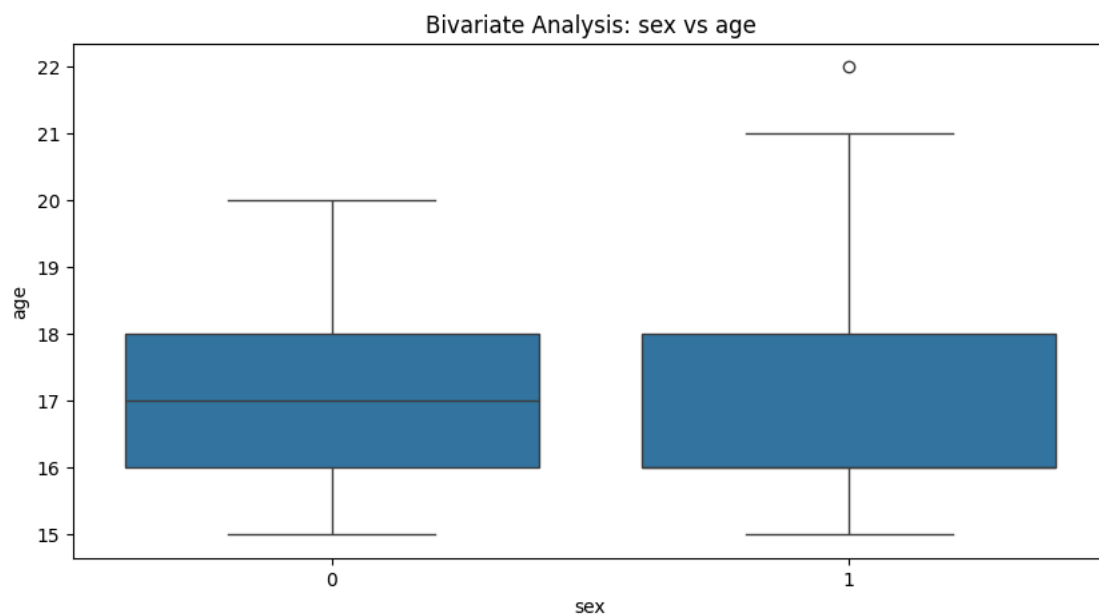


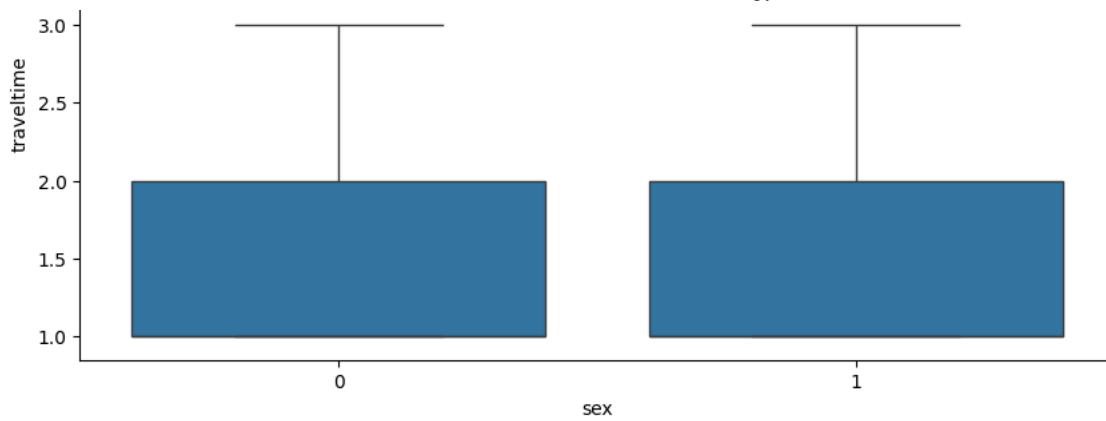
Bivariate Analysis: Age vs G2



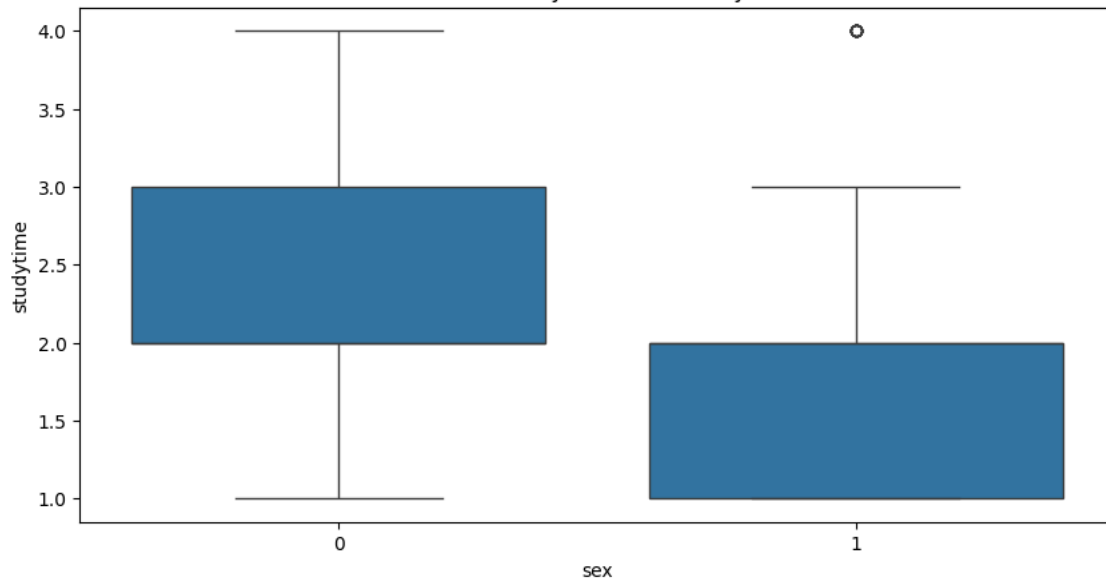


```
categorical_cols = ['sex', 'address', 'famsize', 'schoolsup', 'age_group']
for cat_col in categorical_cols:
    if cat_col in df.columns:
        for num_col in df.select_dtypes(include=['float64', 'int64']).columns:
            plt.figure(figsize=(10,5))
            sns.boxplot(x=df[cat_col].astype(str), y=df[num_col])
            plt.title(f'Bivariate Analysis: {cat_col} vs {num_col}')
            plt.show()
    else:
        print(f"Column '{cat_col}' not found in DataFrame, skipping.")
```

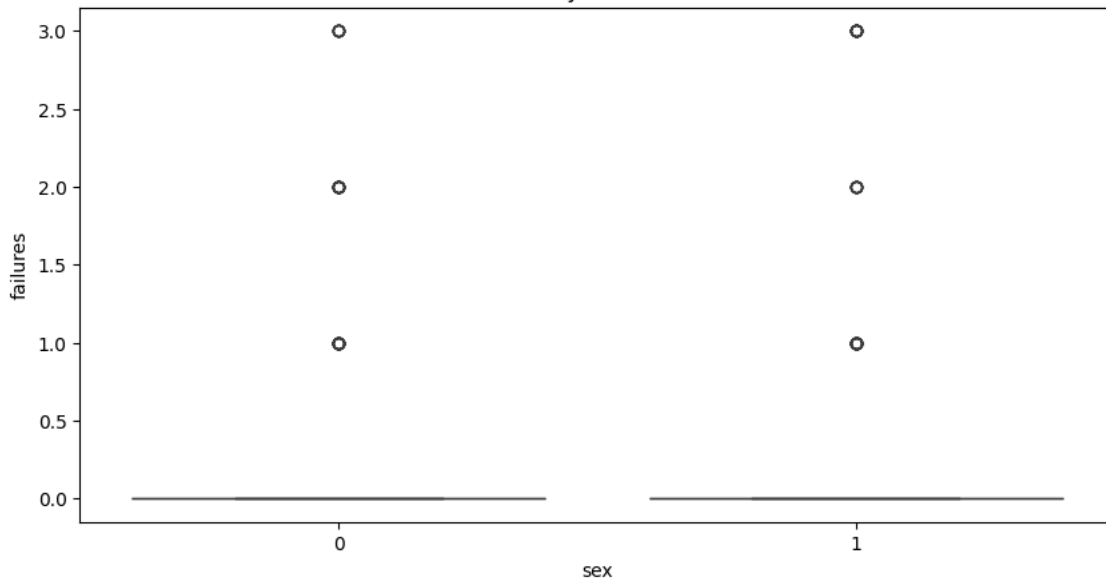




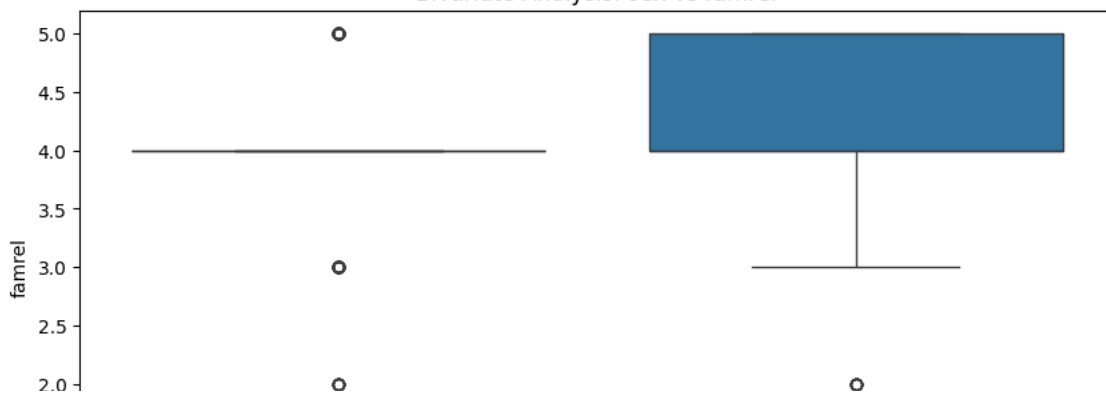
Bivariate Analysis: sex vs studytime

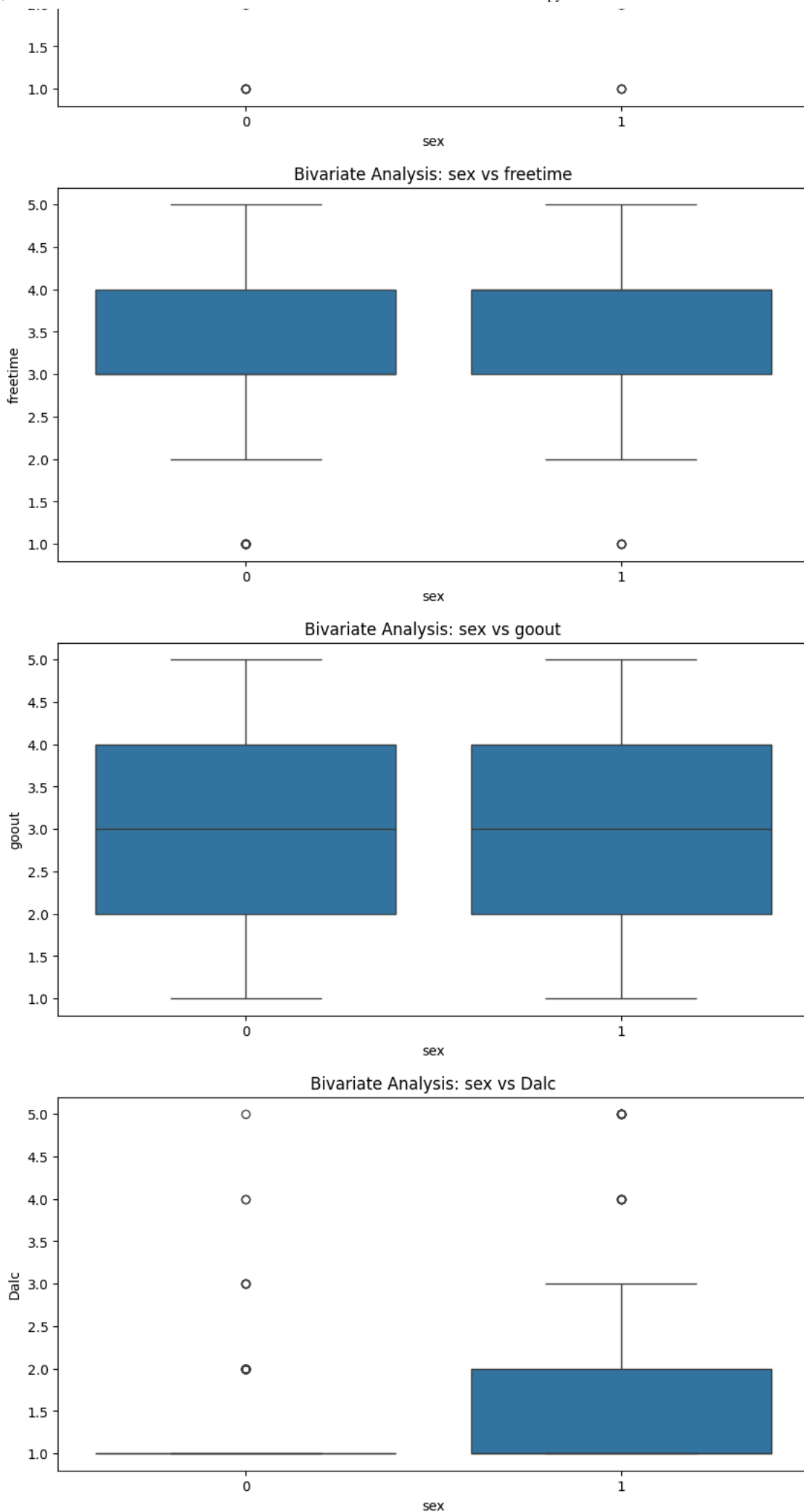


Bivariate Analysis: sex vs failures

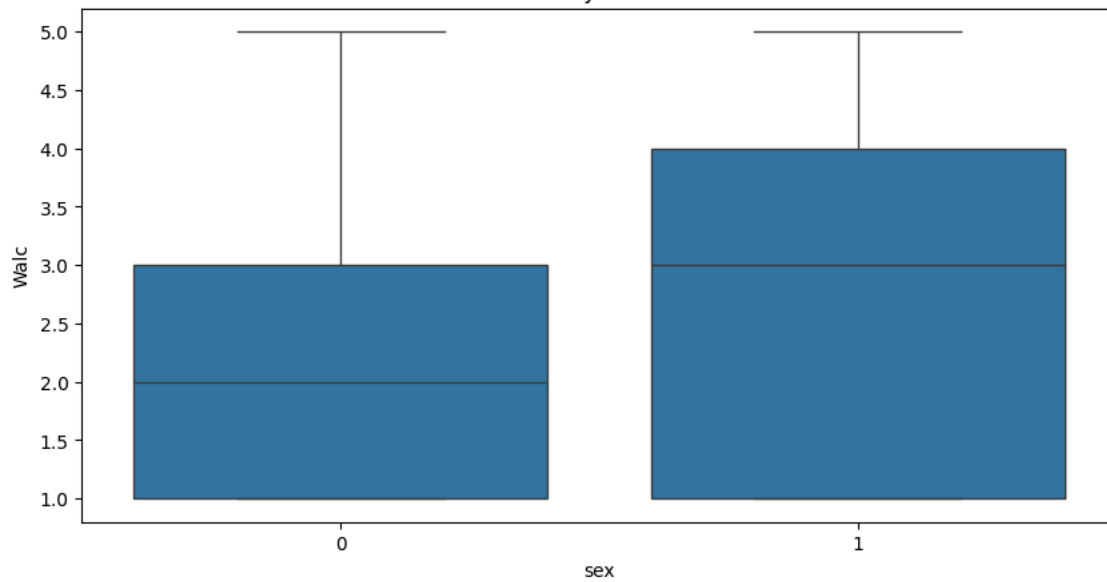


Bivariate Analysis: sex vs famrel

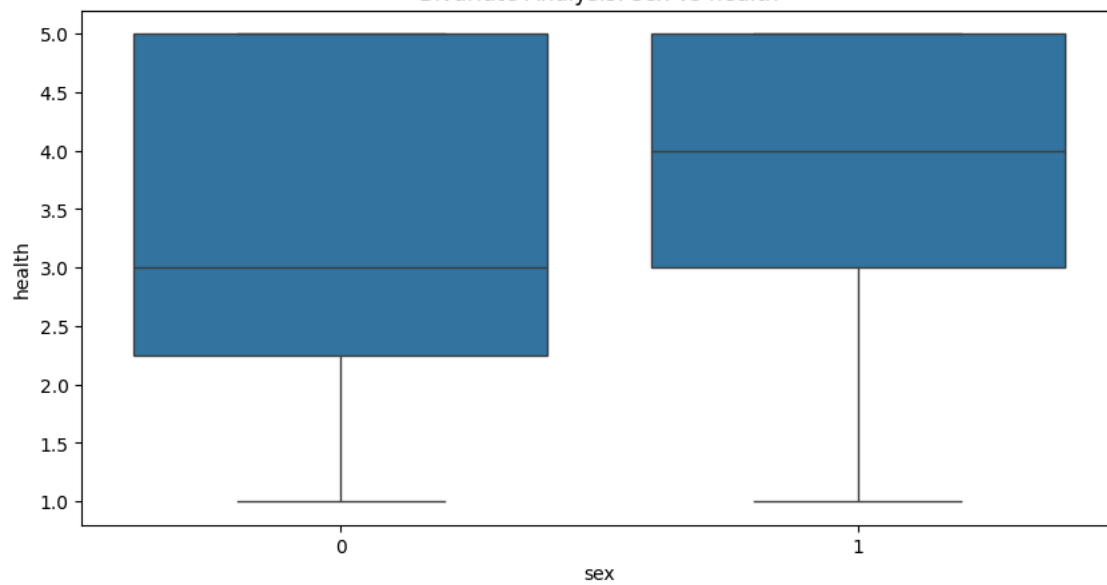




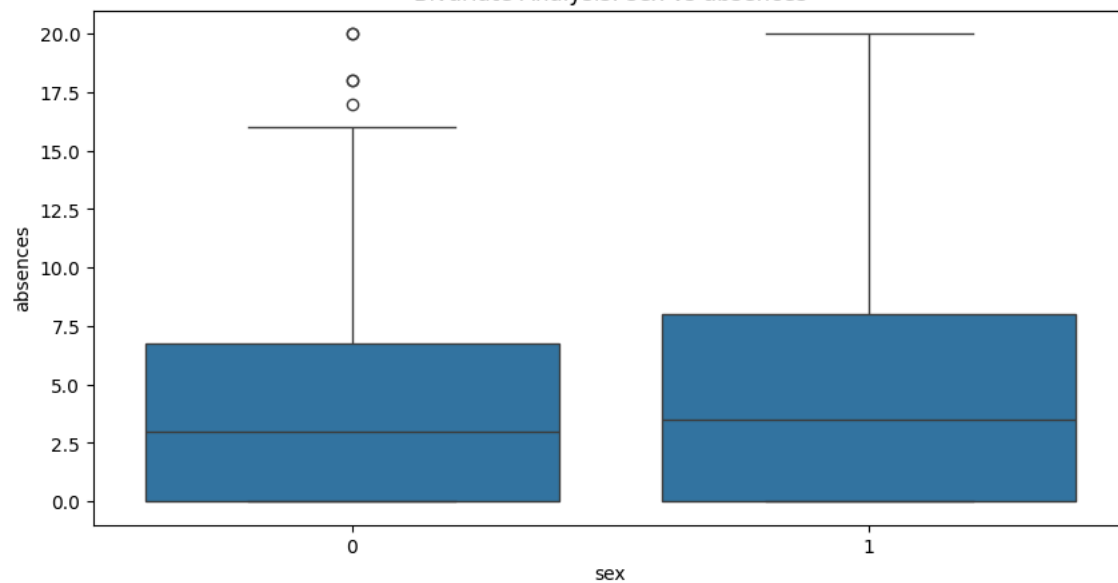
Bivariate Analysis: sex vs Walc



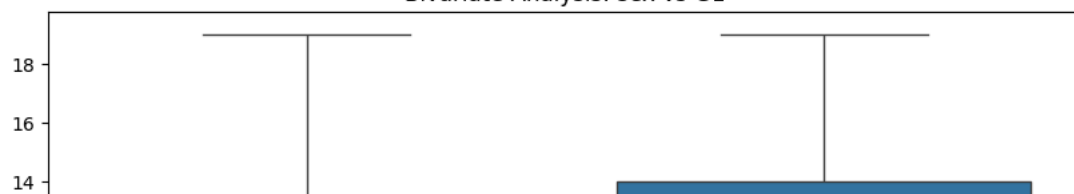
Bivariate Analysis: sex vs health

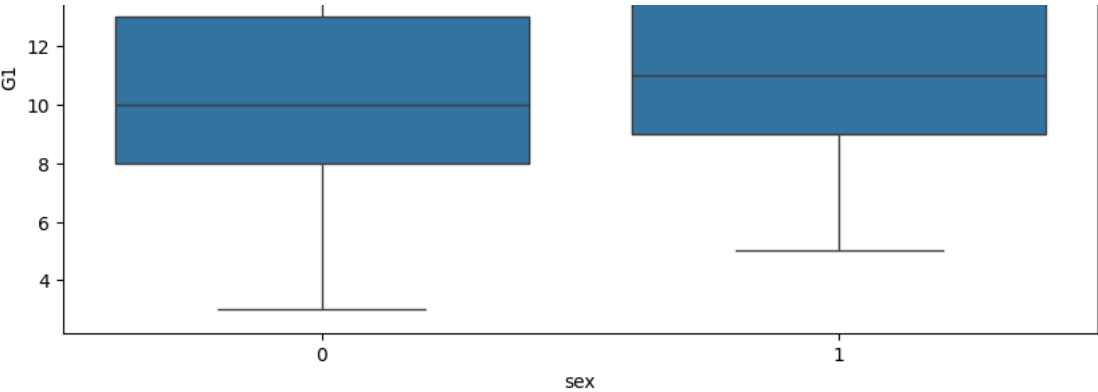


Bivariate Analysis: sex vs absences

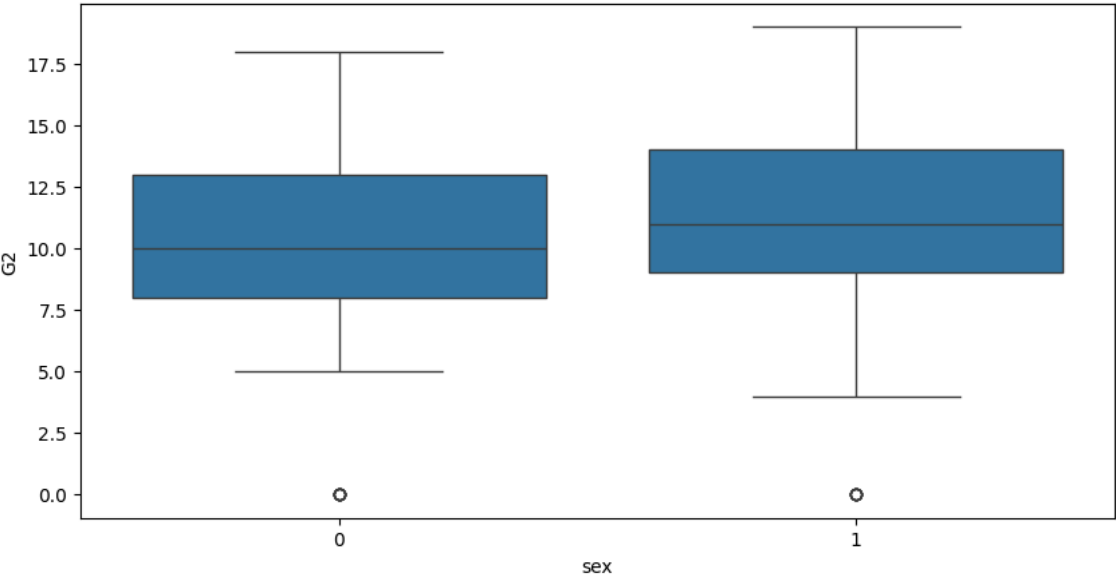


Bivariate Analysis: sex vs G1

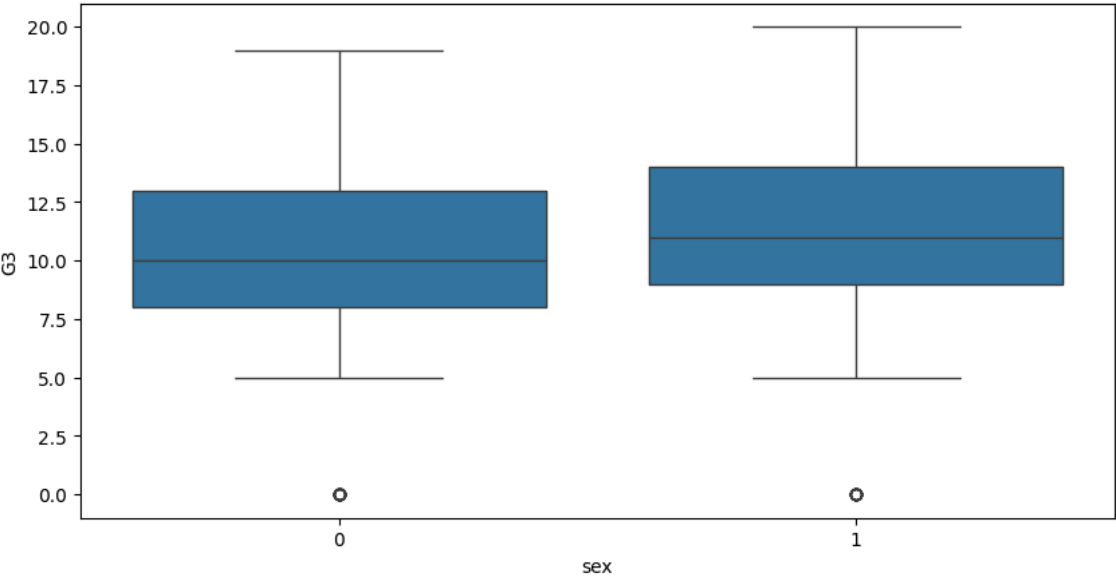




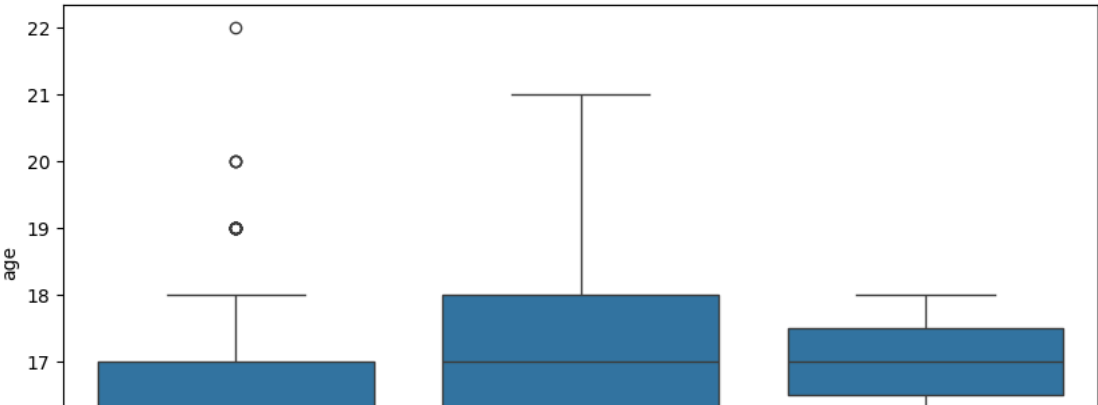
Bivariate Analysis: sex vs G2



Bivariate Analysis: sex vs G3

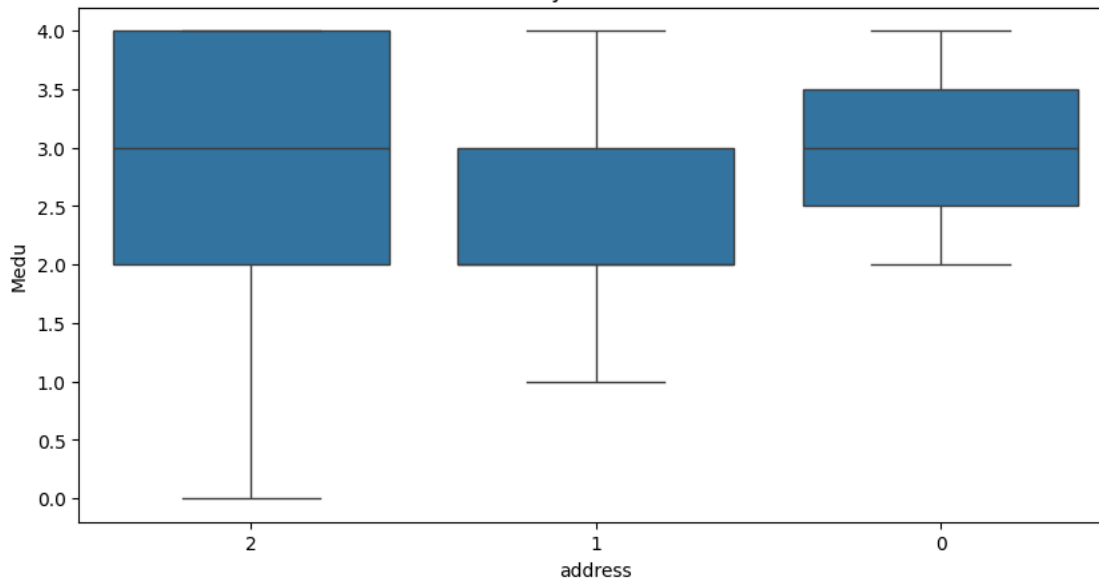


Bivariate Analysis: address vs age

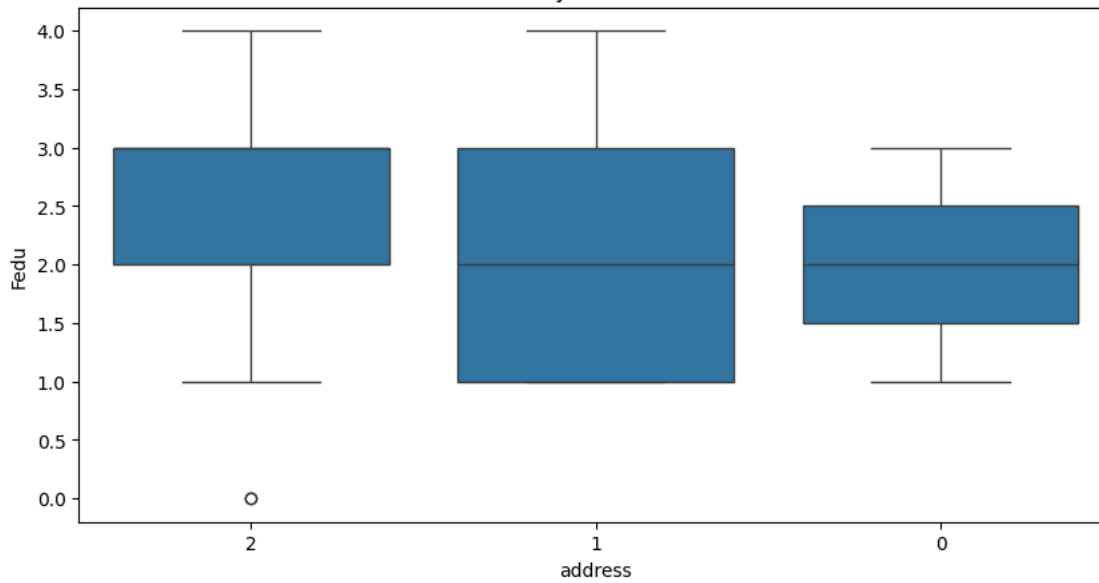




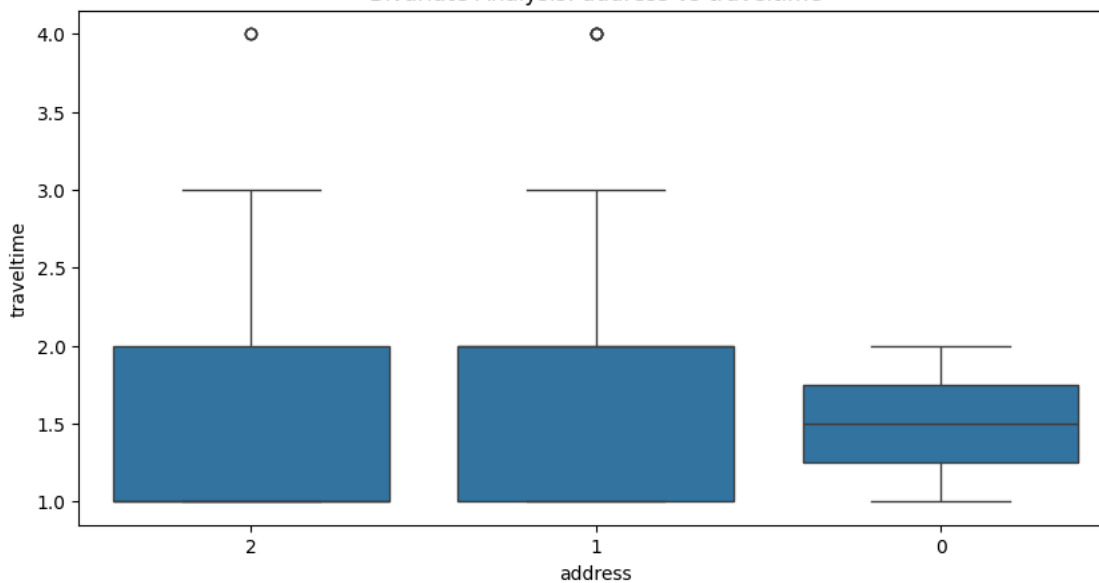
Bivariate Analysis: address vs Medu



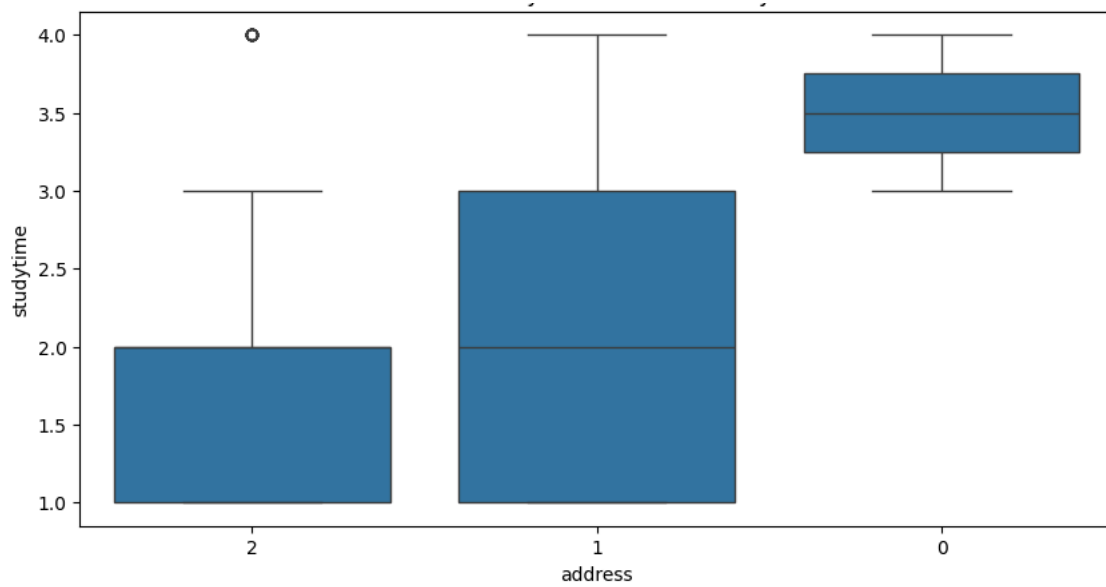
Bivariate Analysis: address vs Fedu



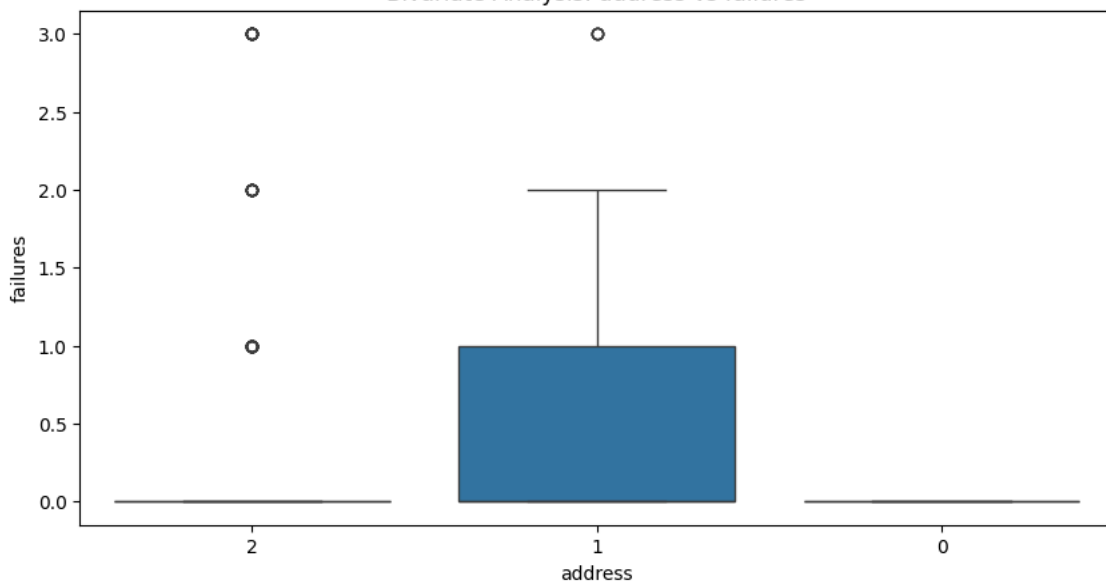
Bivariate Analysis: address vs traveltime



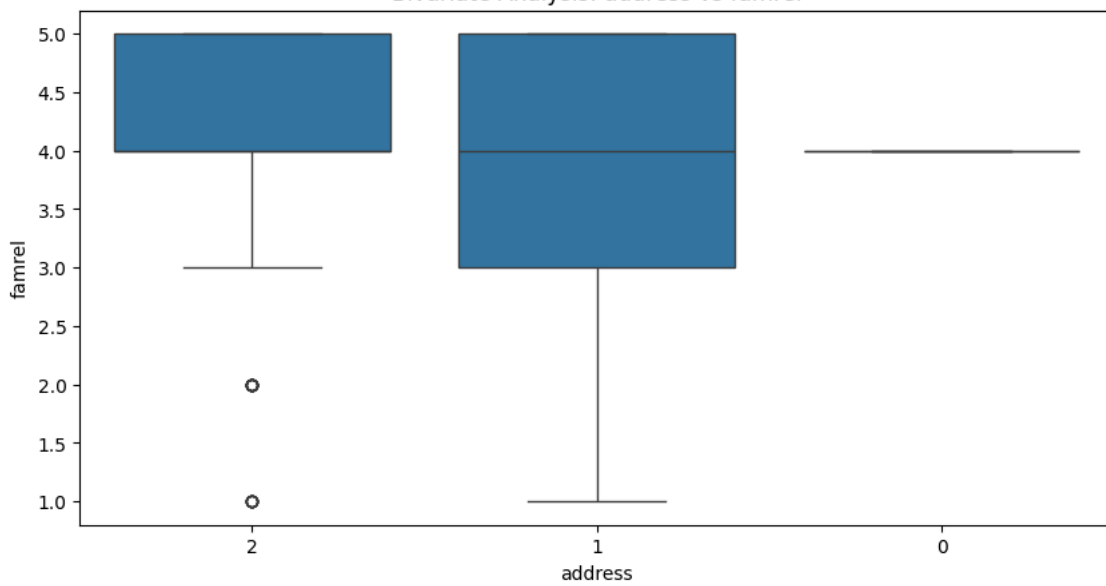
Bivariate Analysis: address vs studytime



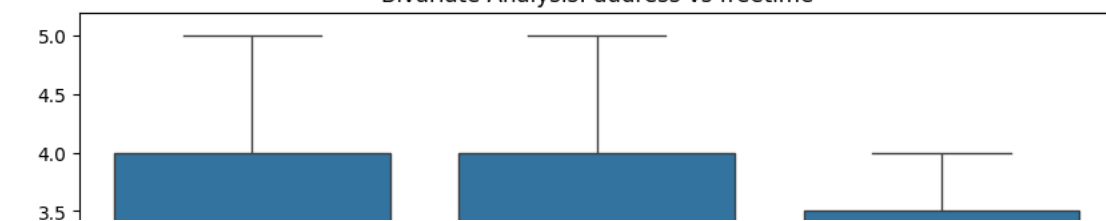
Bivariate Analysis: address vs failures

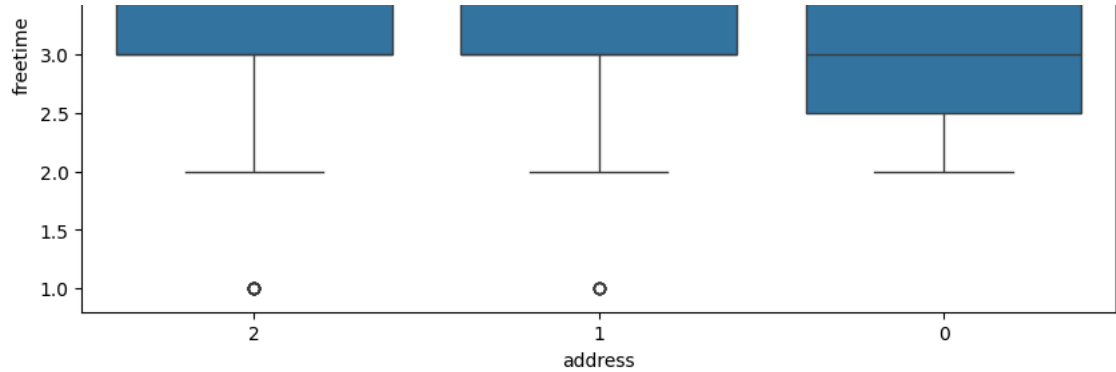


Bivariate Analysis: address vs famrel

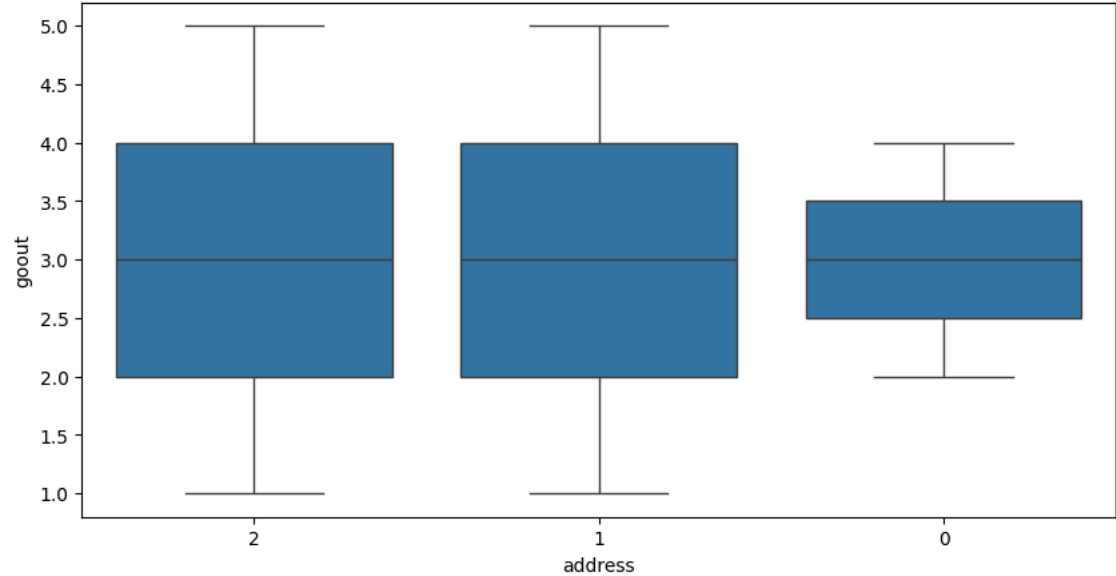


Bivariate Analysis: address vs freetime

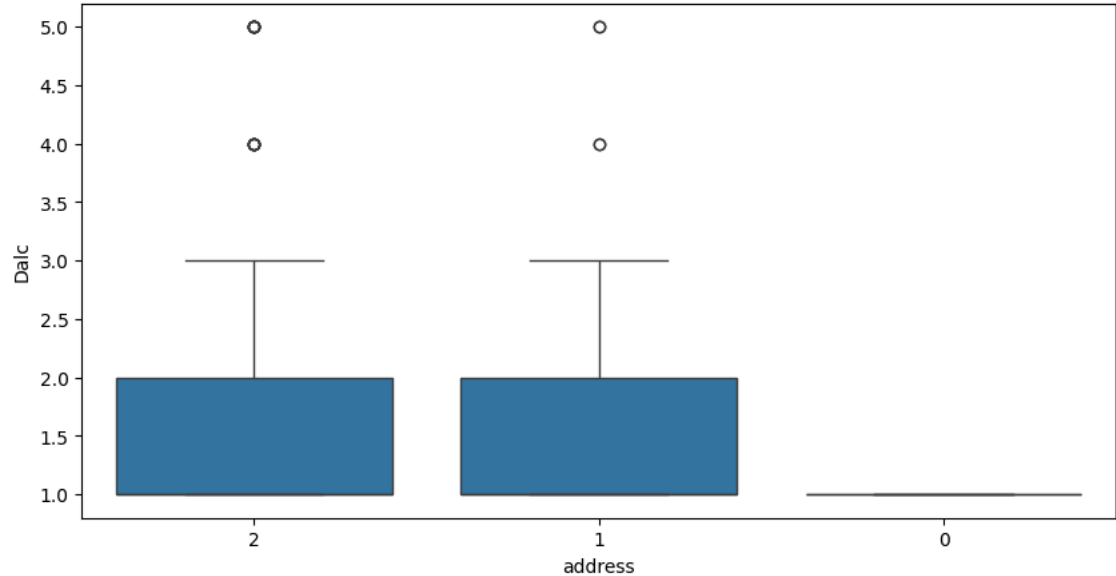




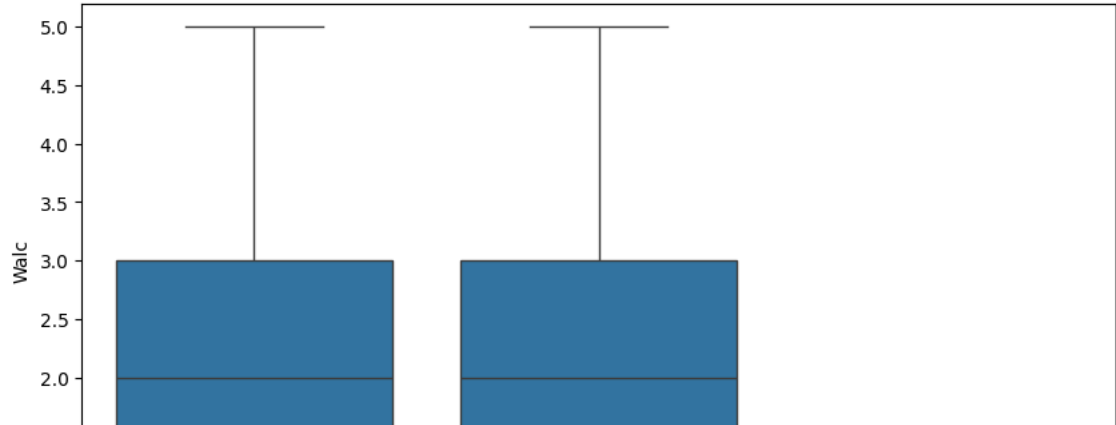
Bivariate Analysis: address vs goout

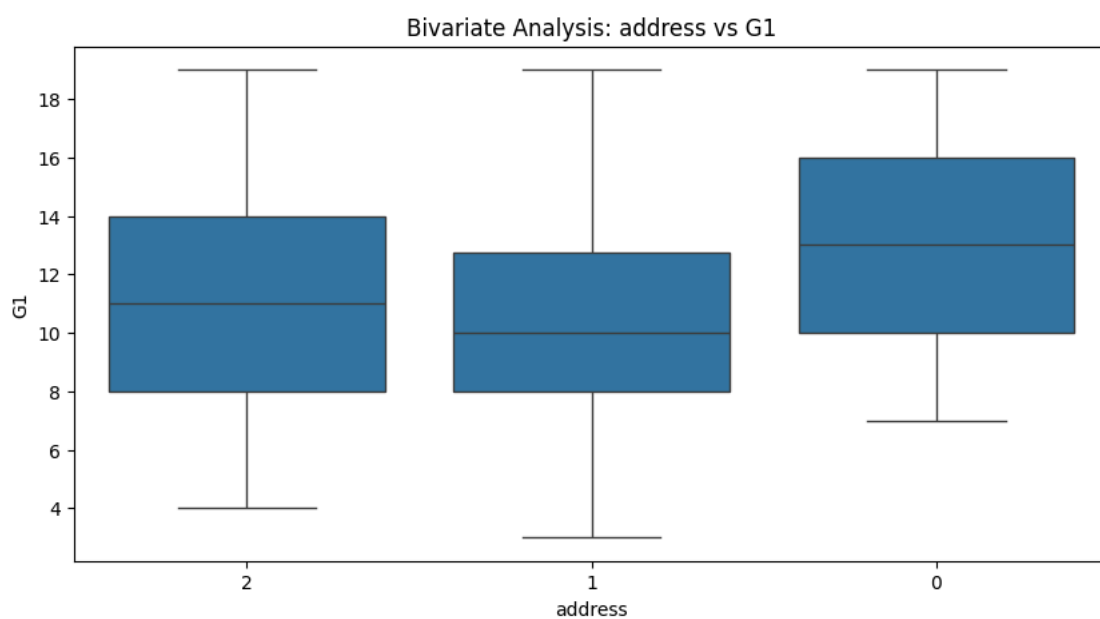
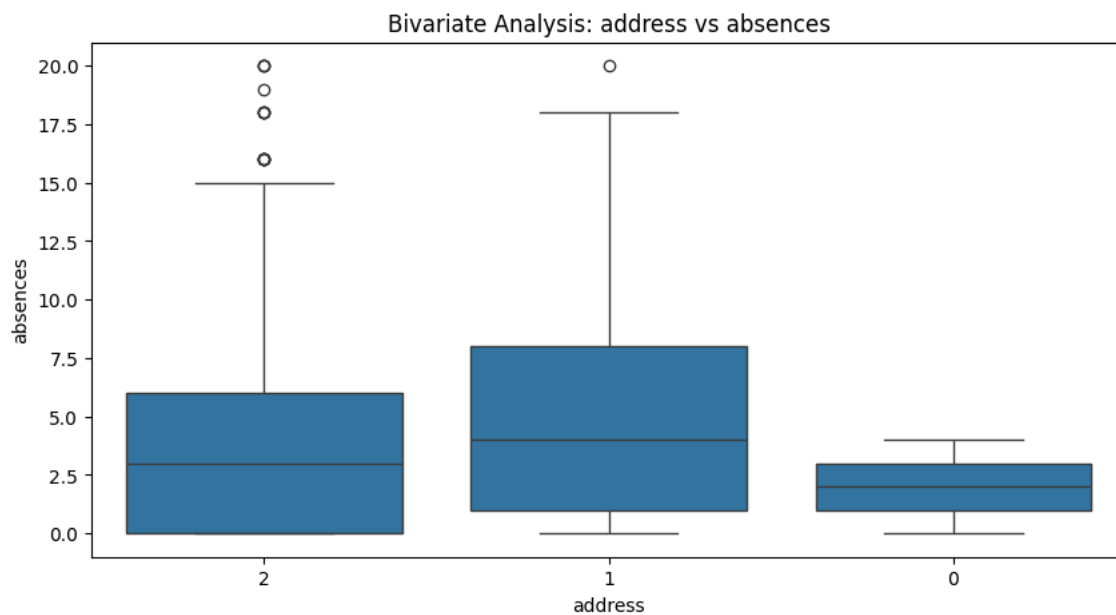
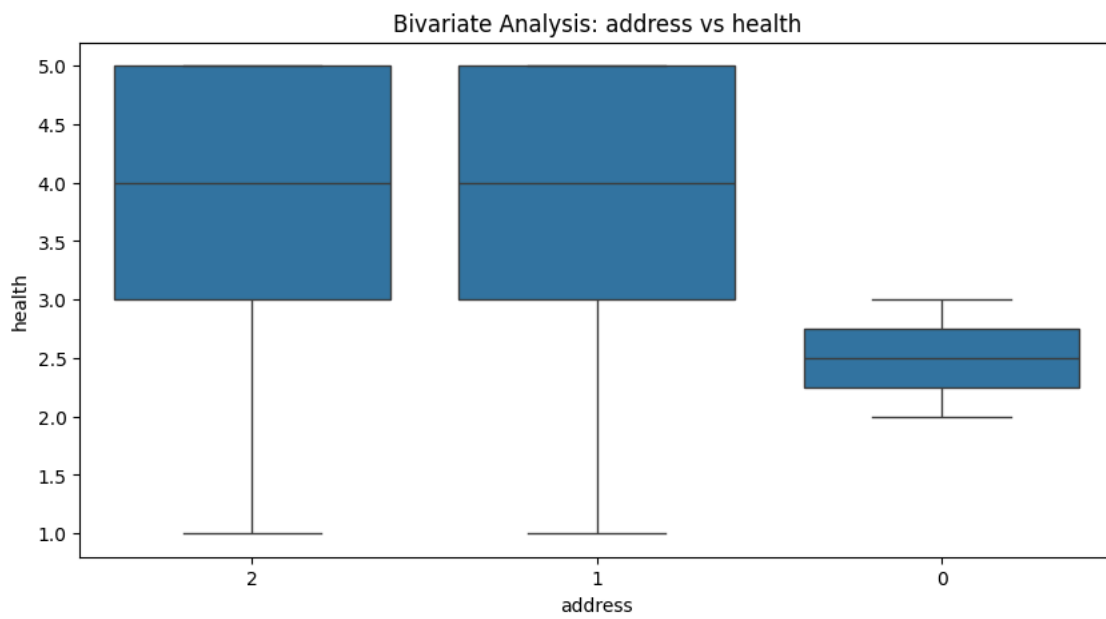


Bivariate Analysis: address vs Dalc

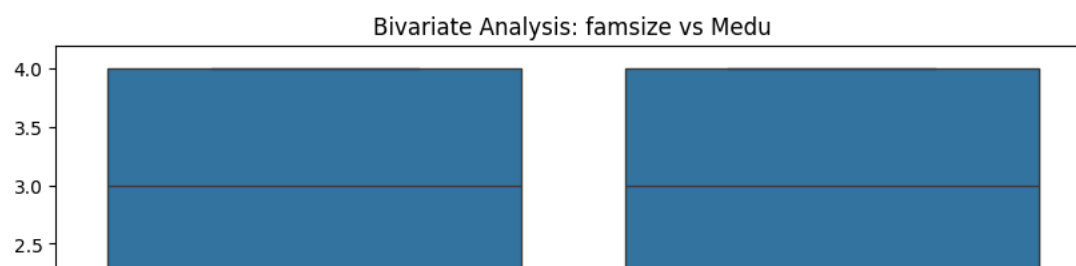
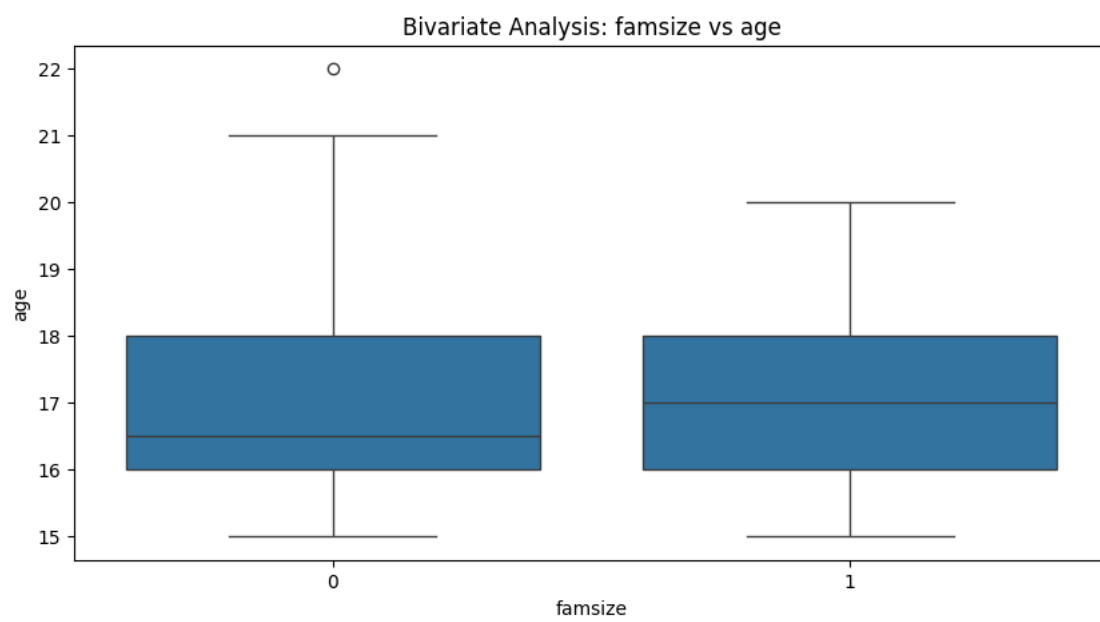
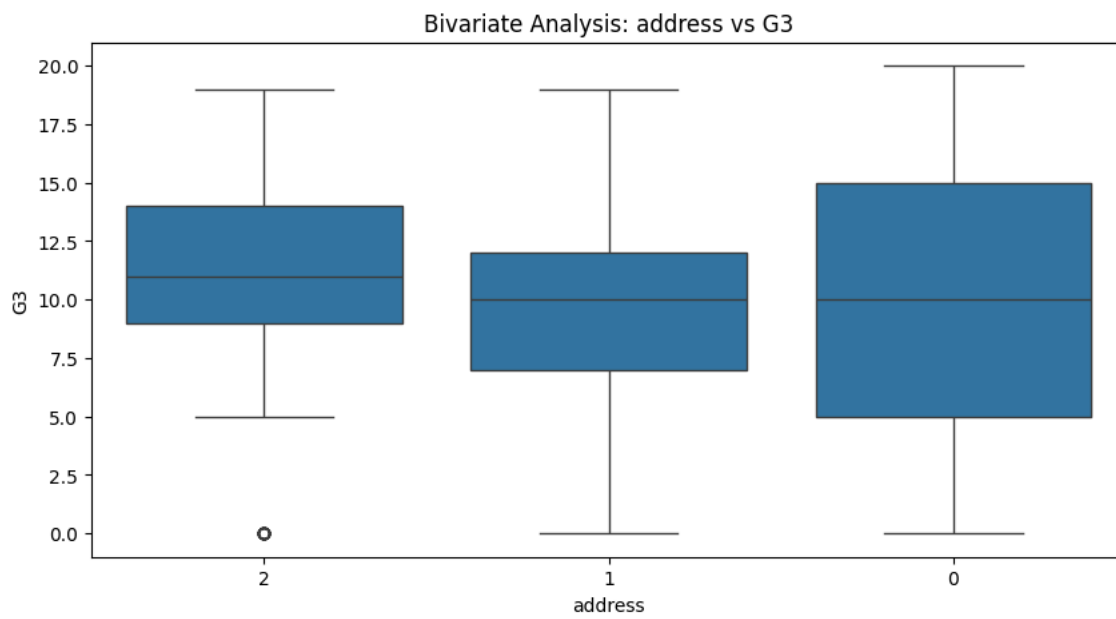
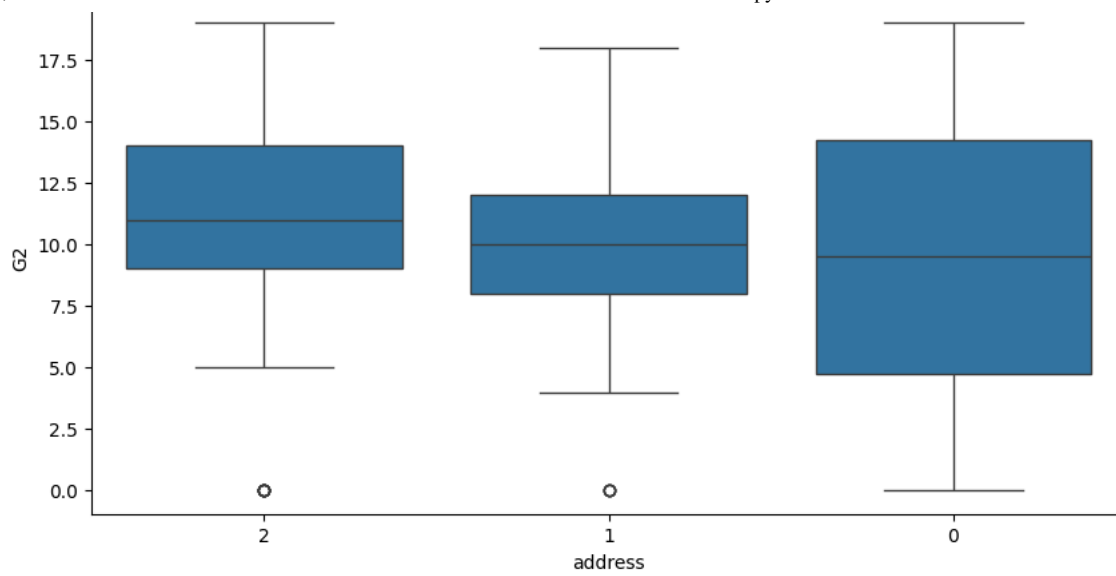


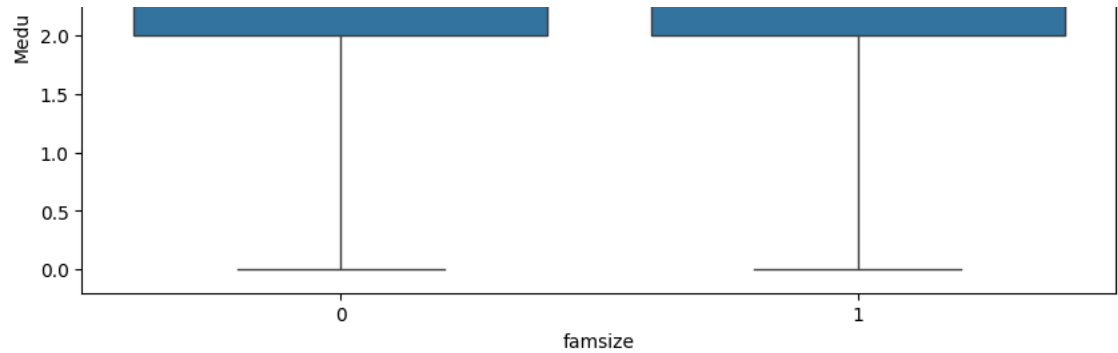
Bivariate Analysis: address vs Walc



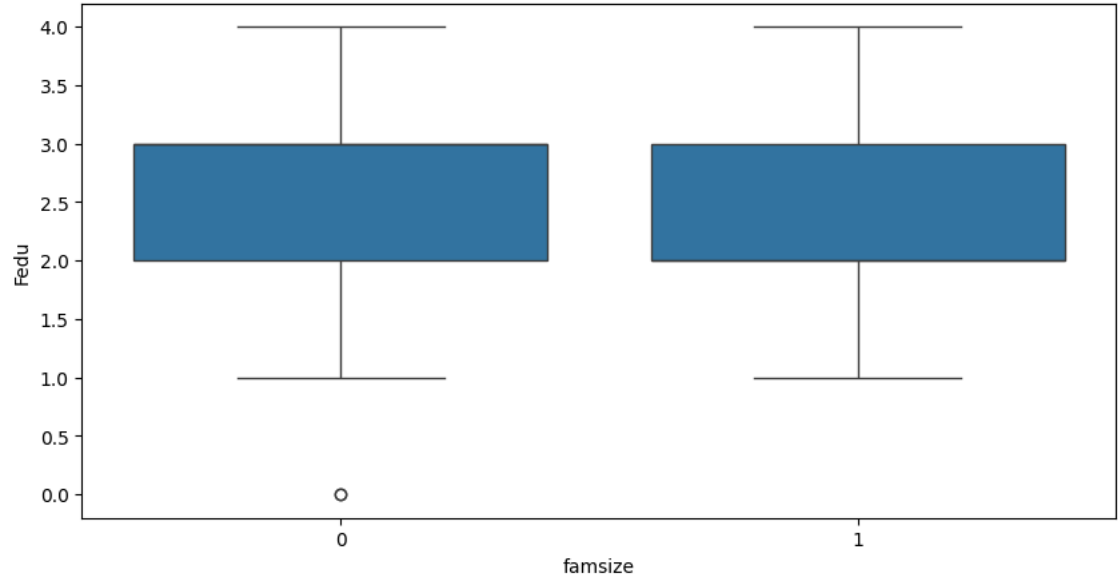


Bivariate Analysis: address vs G2

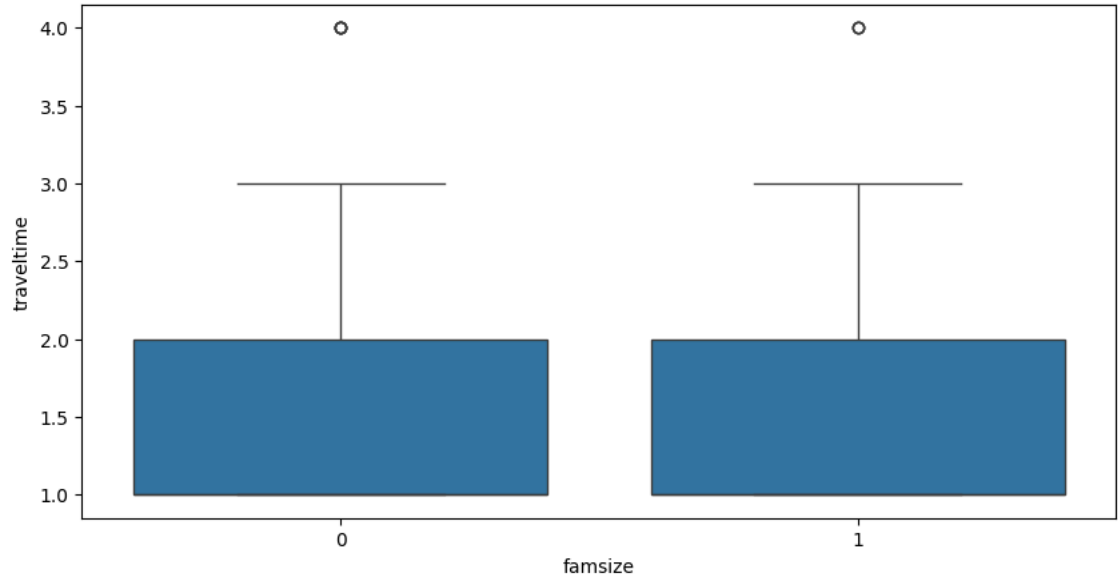




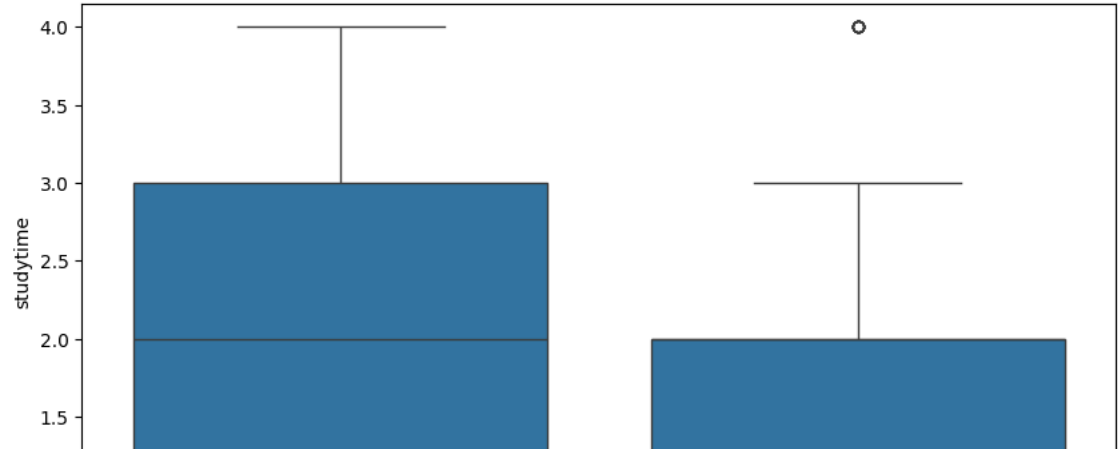
Bivariate Analysis: famsize vs Fedu

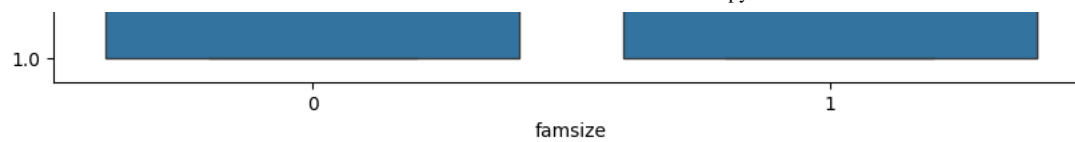


Bivariate Analysis: famsize vs traveltime

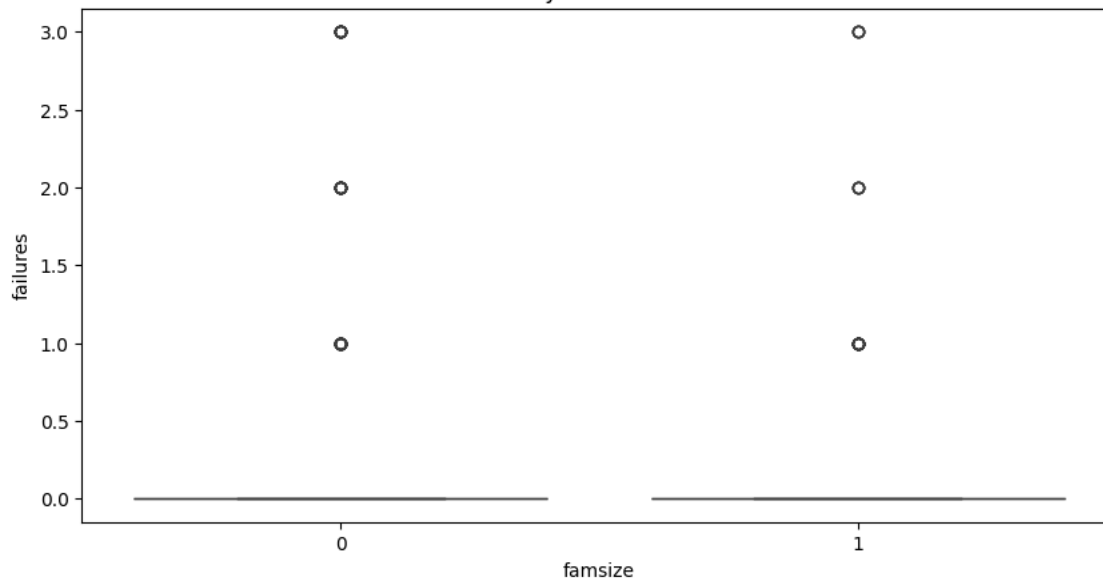


Bivariate Analysis: famsize vs studytime

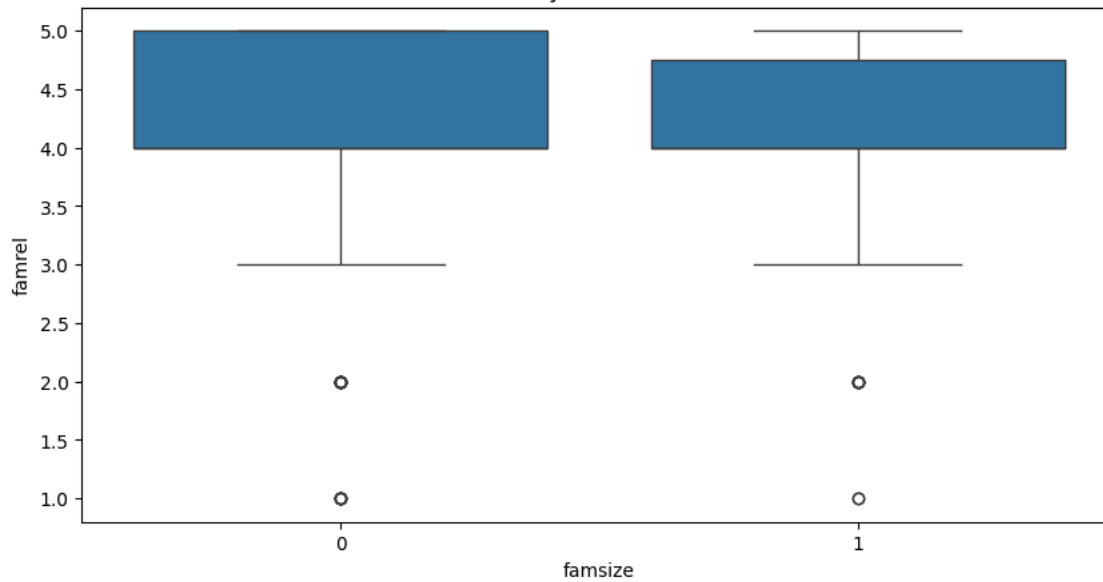




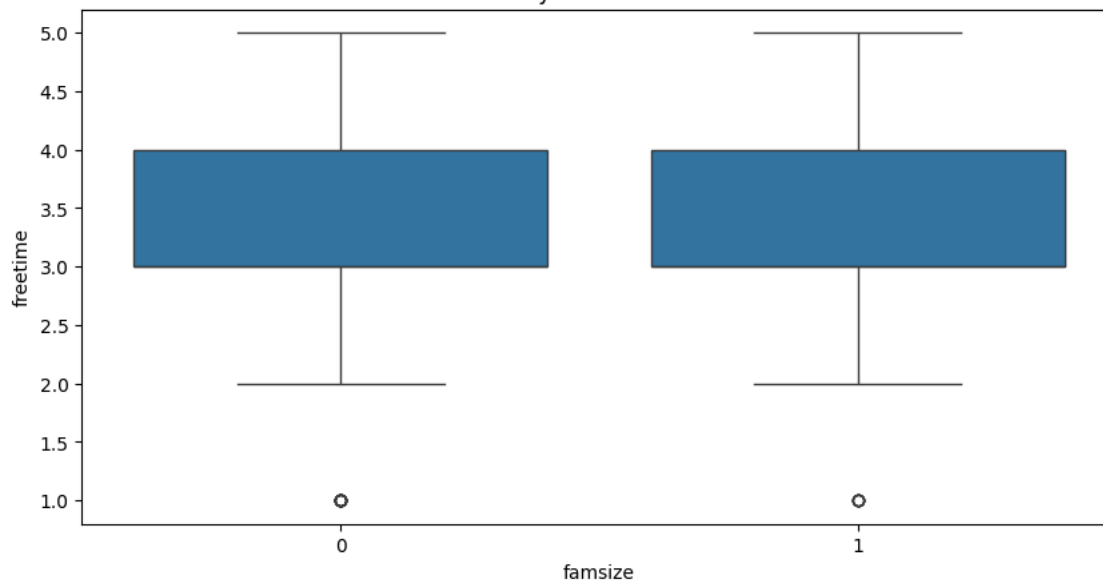
Bivariate Analysis: famsize vs failures



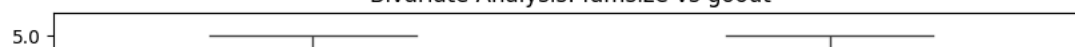
Bivariate Analysis: famsize vs famrel

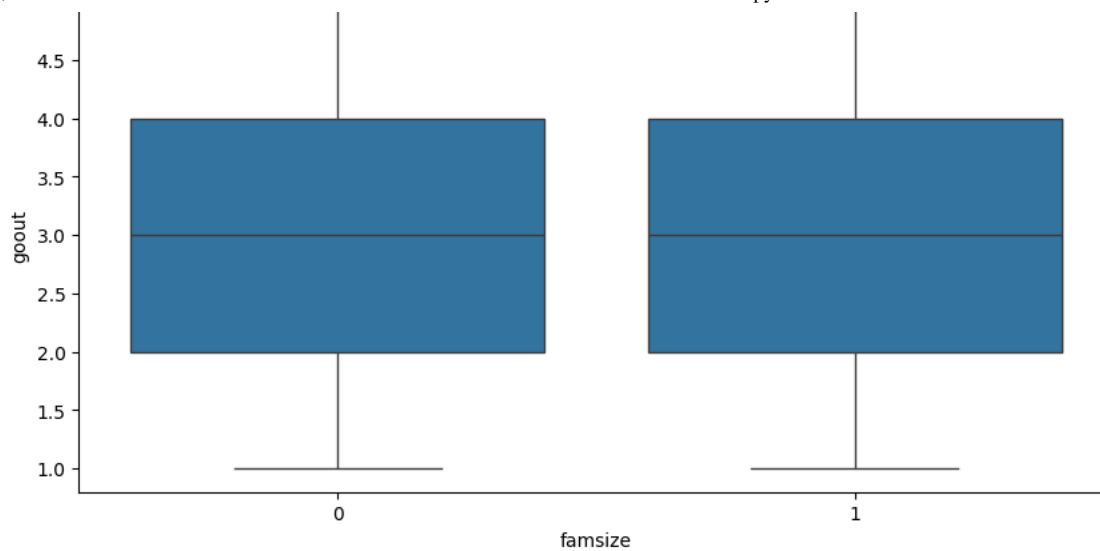


Bivariate Analysis: famsize vs freetime

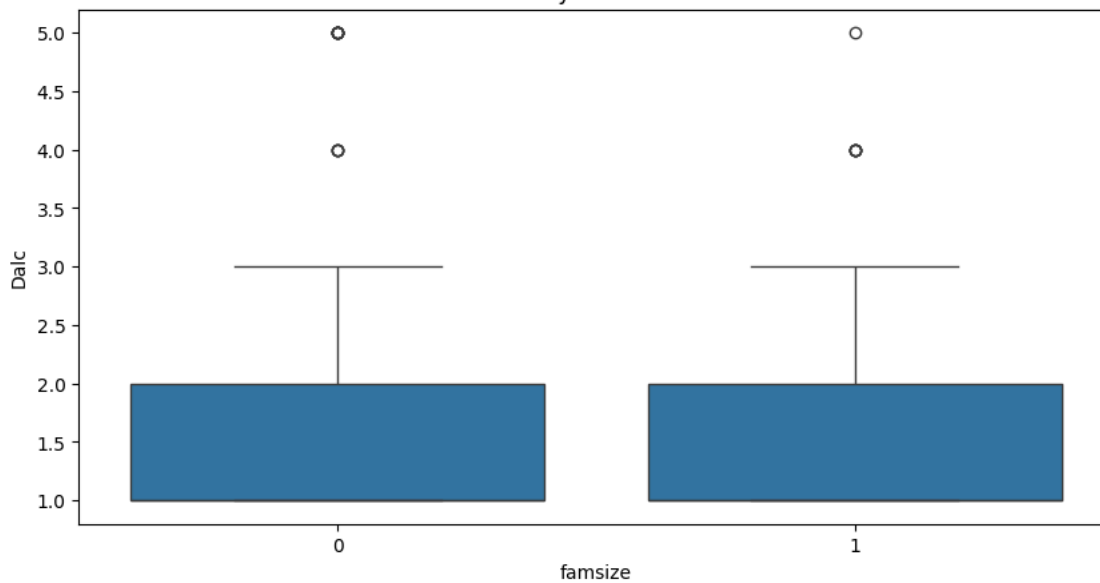


Bivariate Analysis: famsize vs goout

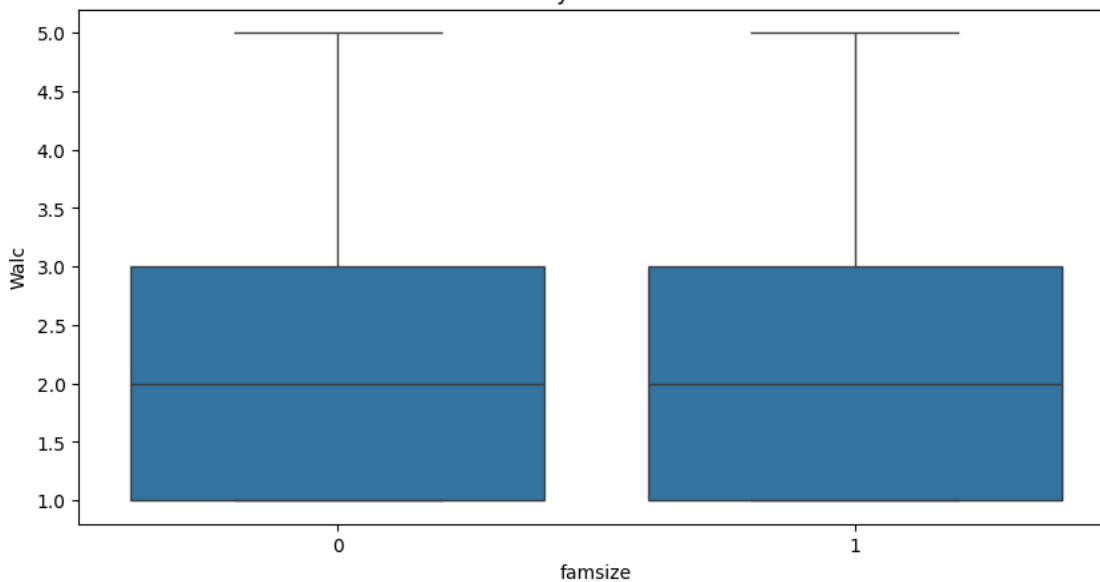




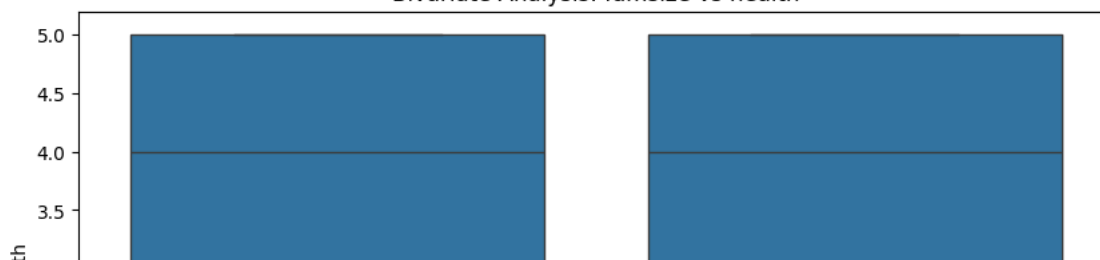
Bivariate Analysis: famsize vs Dalc

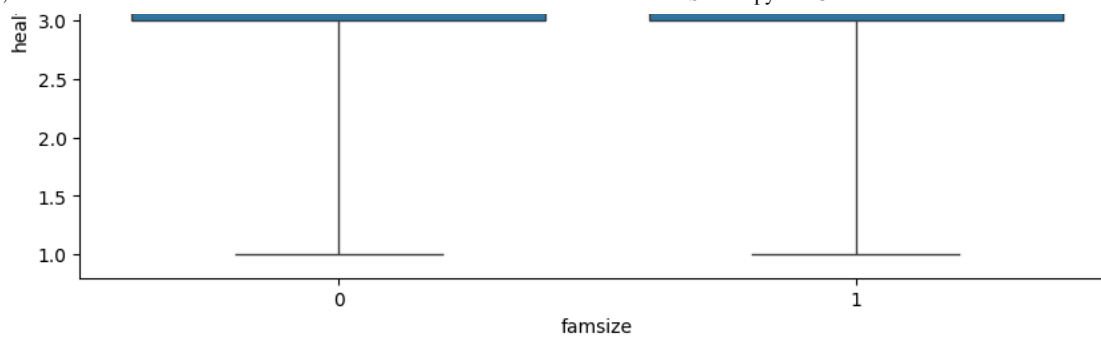


Bivariate Analysis: famsize vs Walc

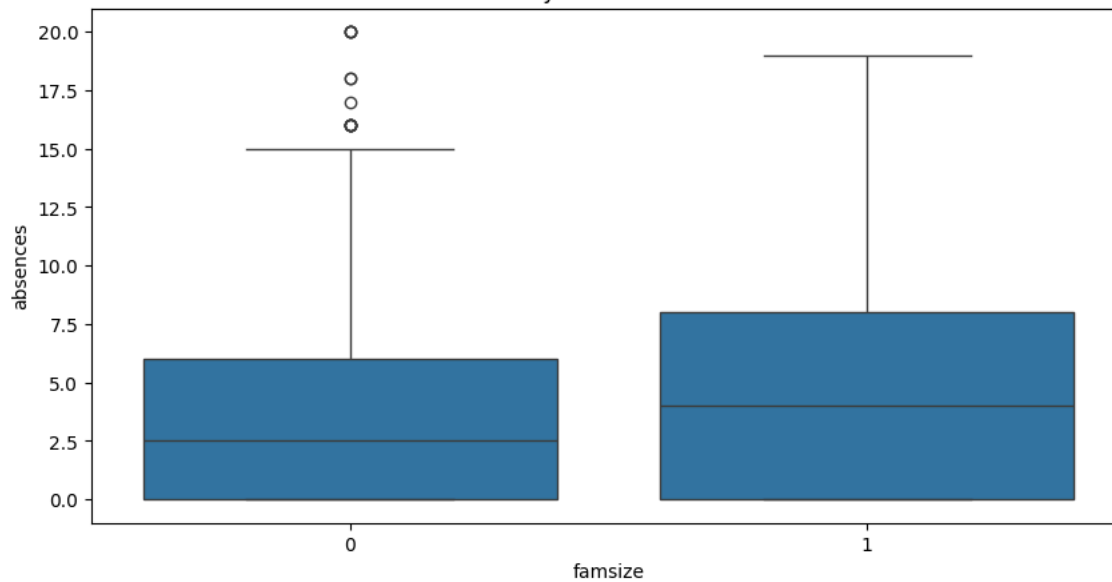


Bivariate Analysis: famsize vs health

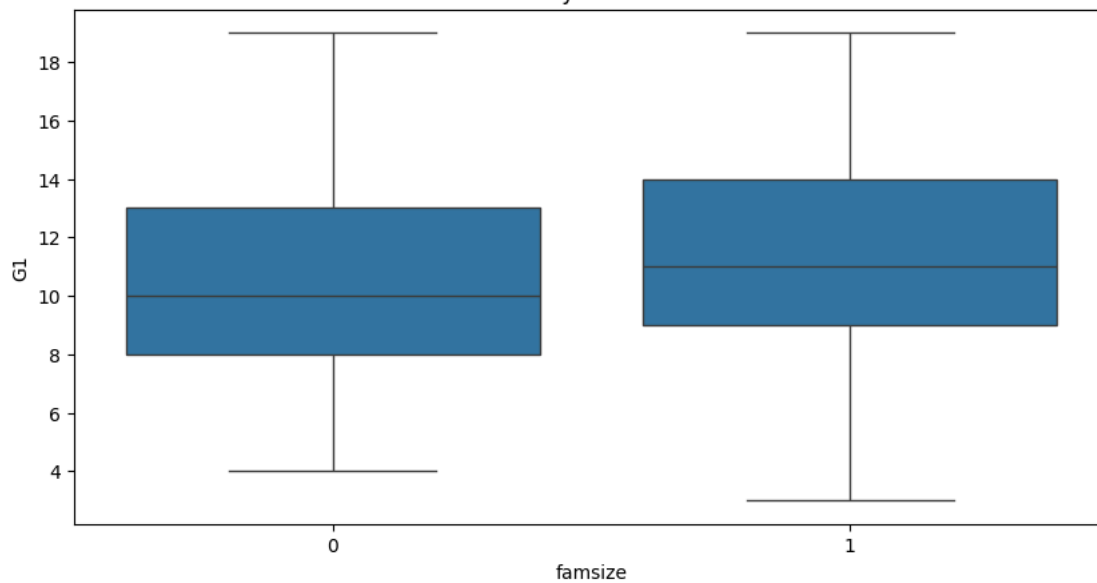




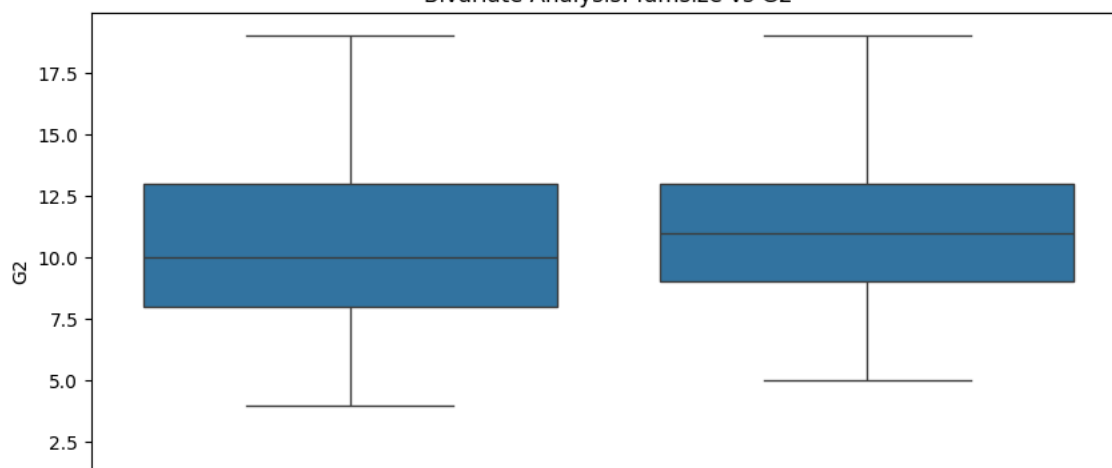
Bivariate Analysis: famsize vs absences

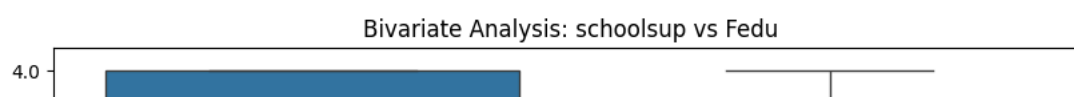
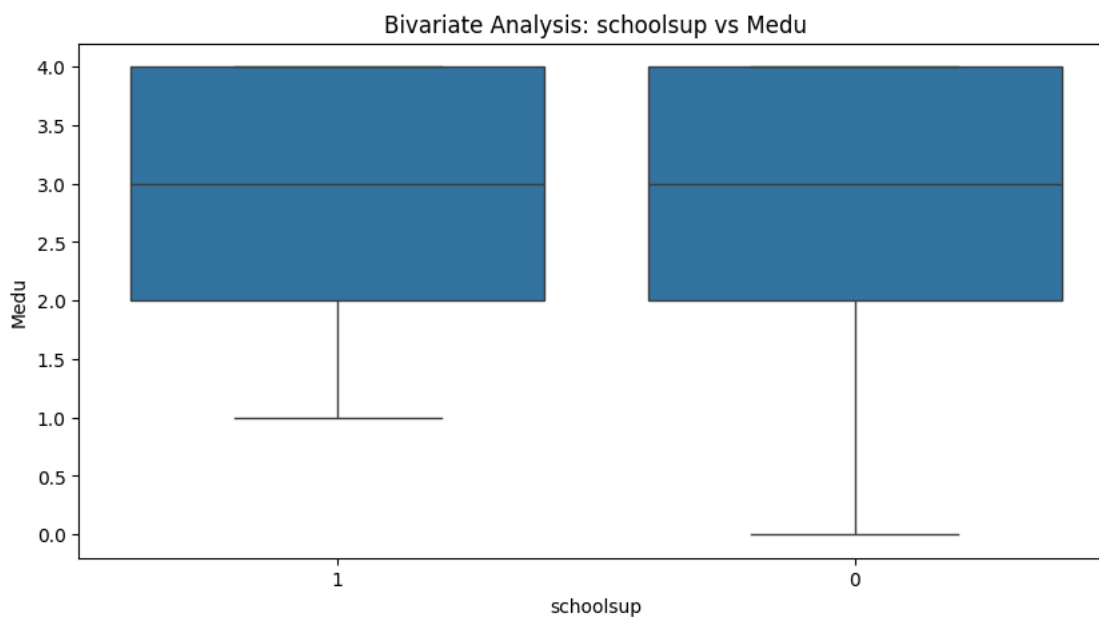
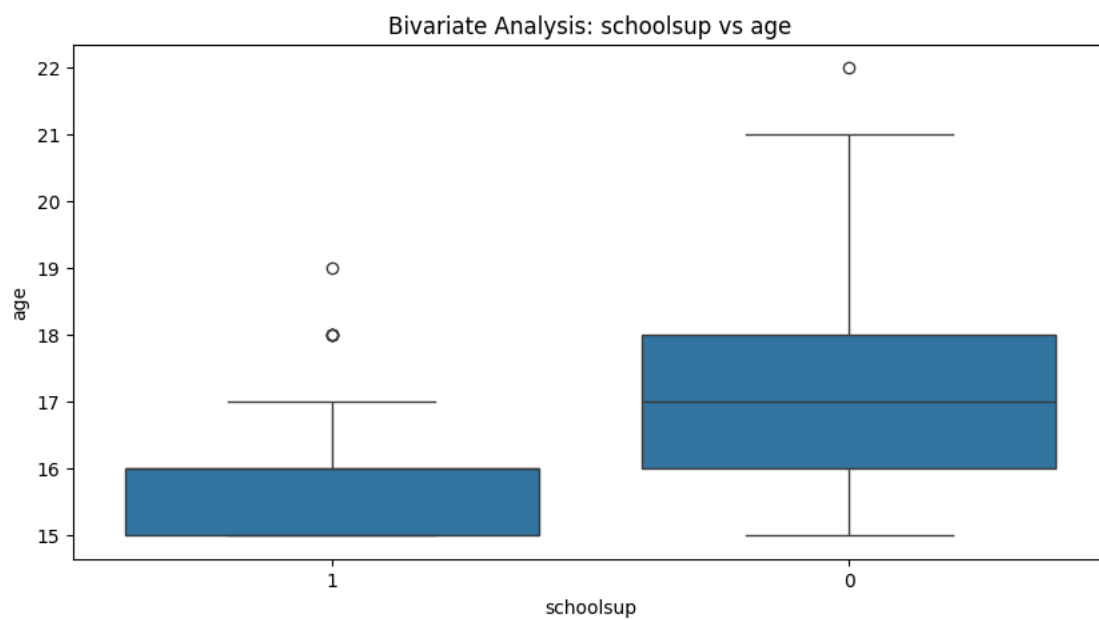
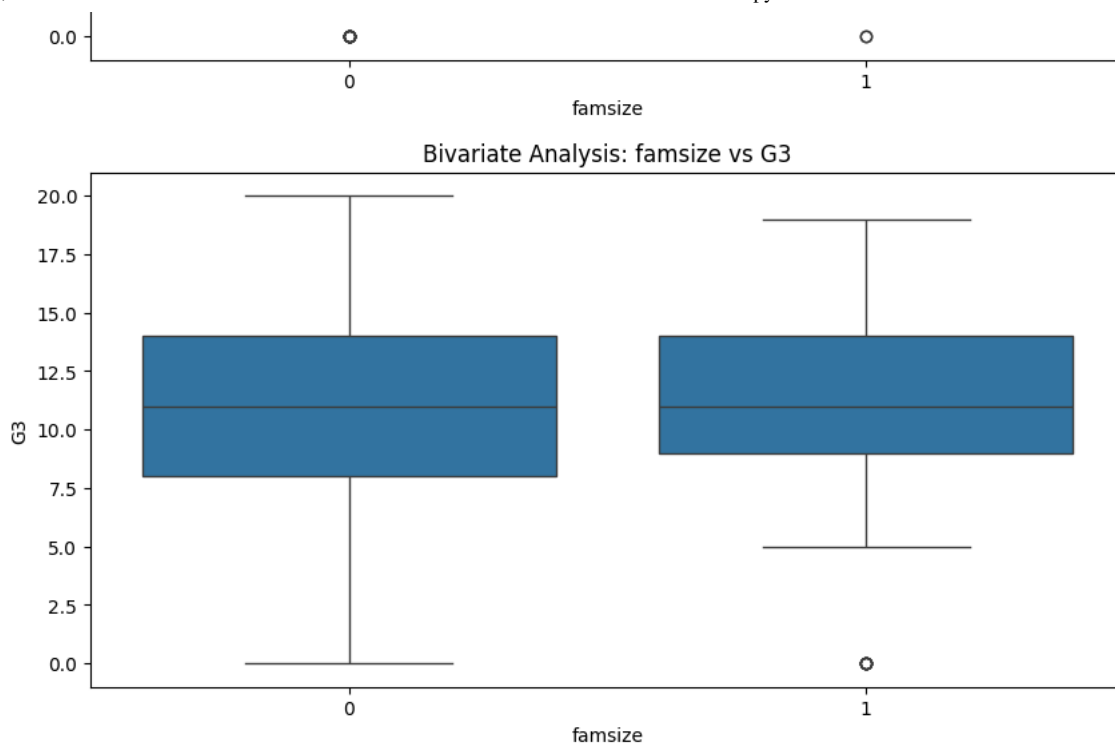


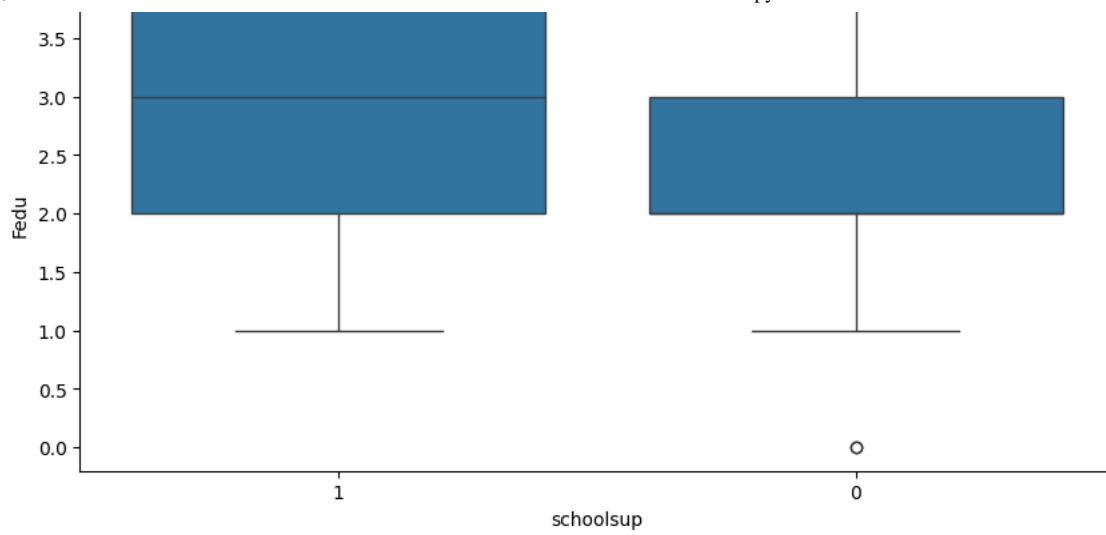
Bivariate Analysis: famsize vs G1



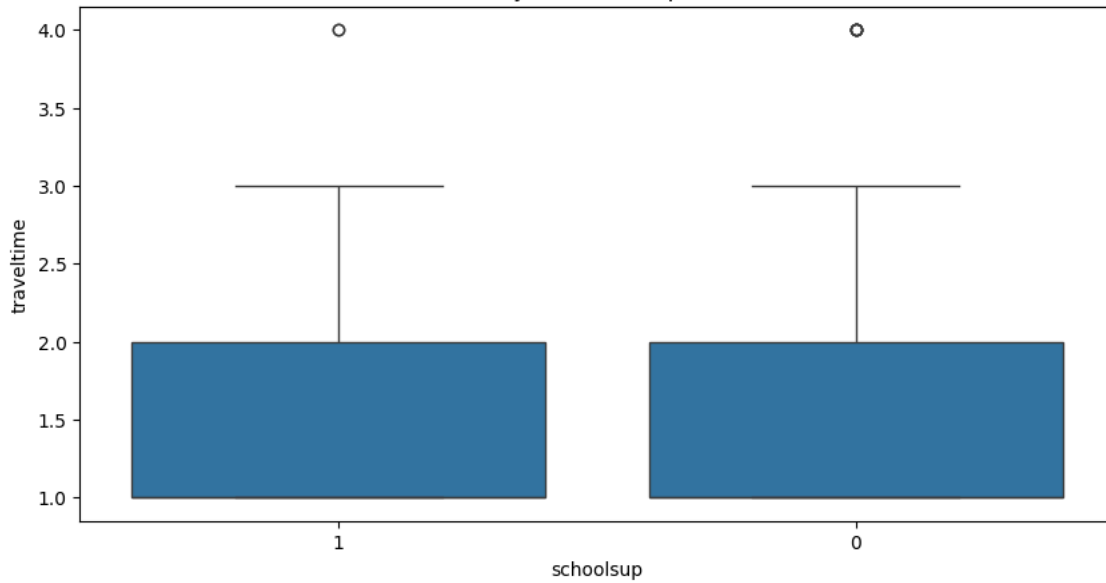
Bivariate Analysis: famsize vs G2



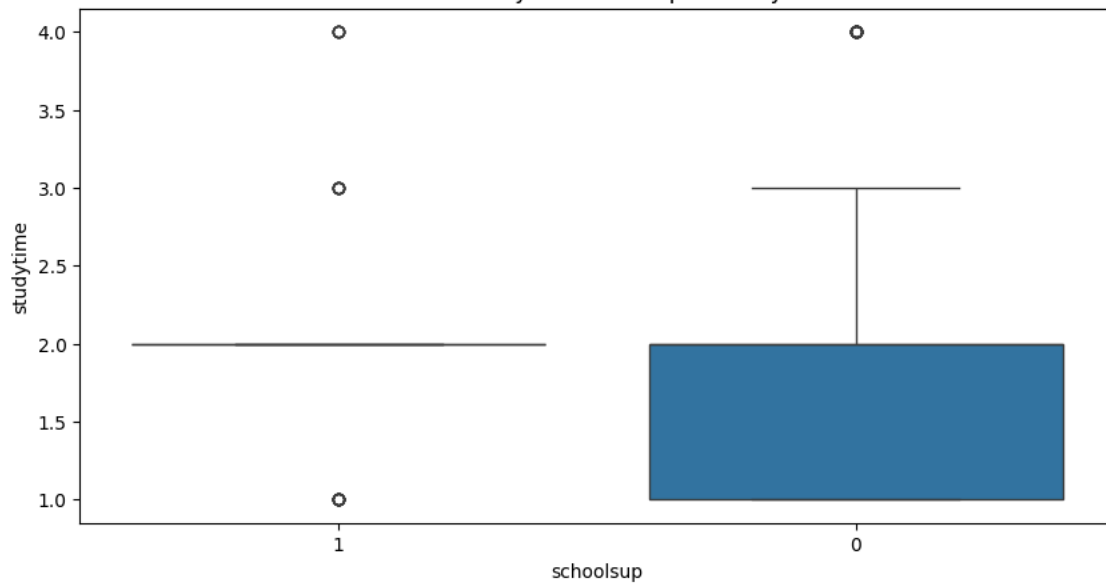




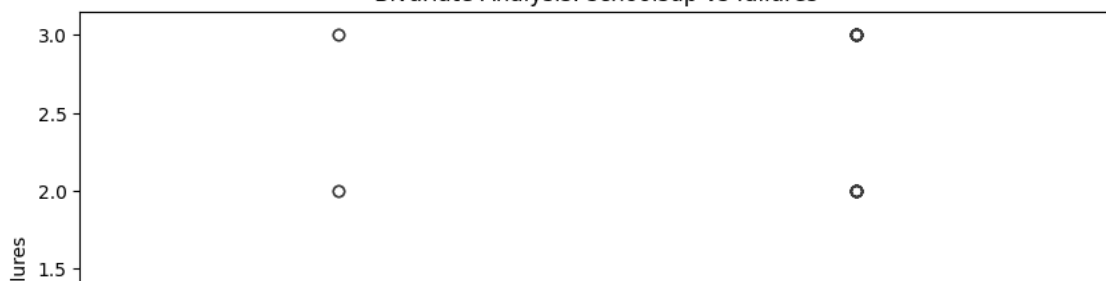
Bivariate Analysis: schoolsup vs traveltime

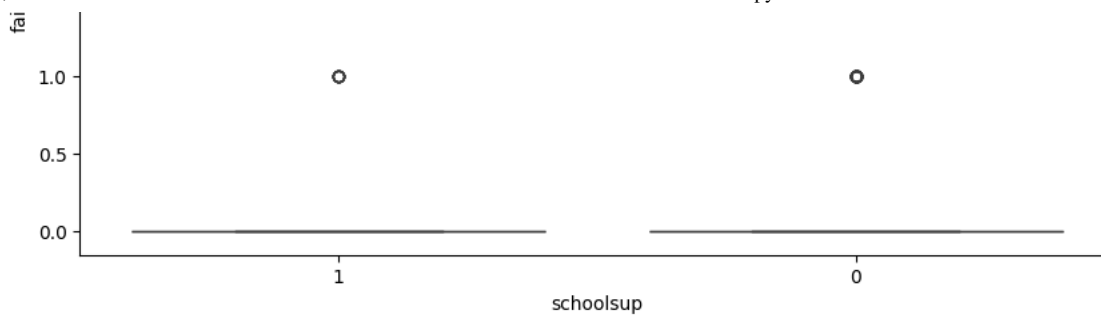


Bivariate Analysis: schoolsup vs studytime

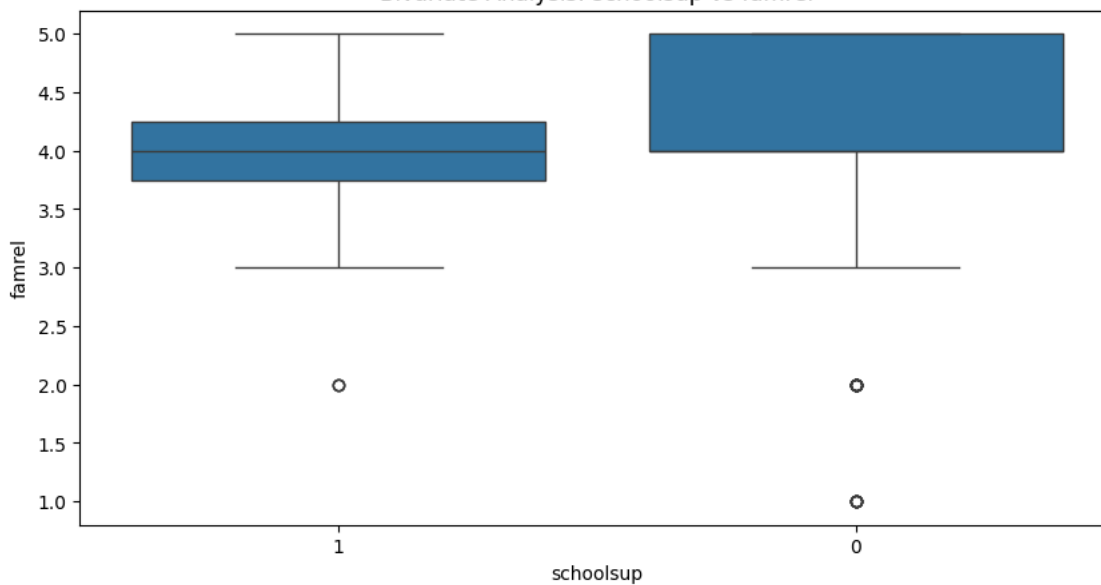


Bivariate Analysis: schoolsup vs failures

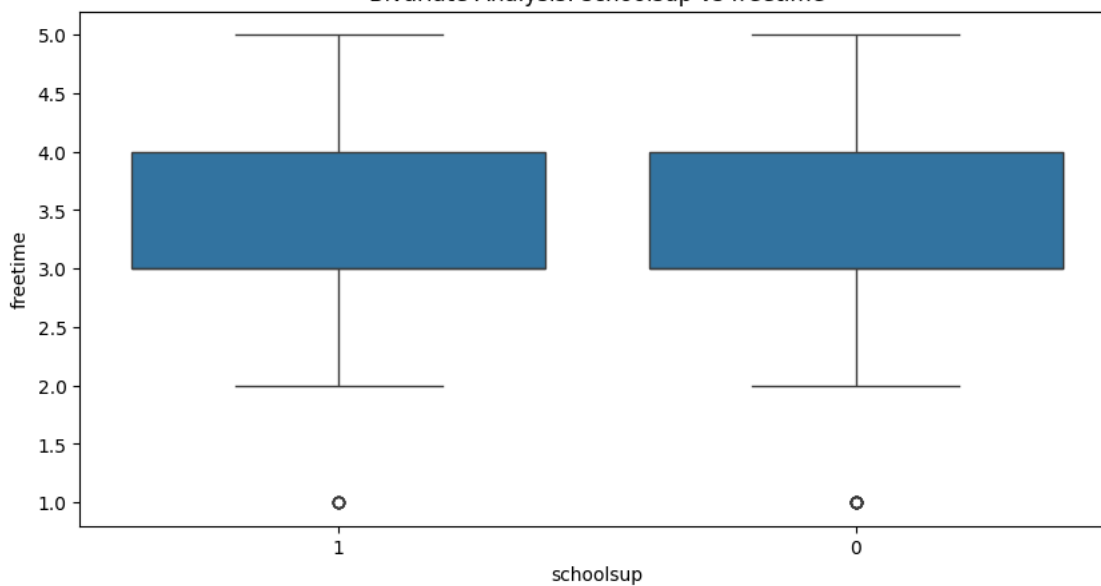




Bivariate Analysis: schoolsup vs famrel



Bivariate Analysis: schoolsup vs freetime



Bivariate Analysis: schoolsup vs goout

