

```
In [1]: import pandas as pd
import numpy as np
import scipy as sp
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: from google.colab import files

uploaded = files.upload()

for fn in uploaded.keys():
    print('User uploaded file "{name}" with length {length} bytes'.format(
        name=fn, length=len(uploaded[fn])))
```

Choose Files No file selected

Upload widget is only available when the cell has

been executed in the current browser session. Please rerun this cell to enable.

Saving Test_scores_Math Scores.xlsx to Test_scores_Math Scores.xlsx

User uploaded file "Test_scores_Math Scores.xlsx" with length 26762 bytes

```
In [13]: # Determine bin edges based on quantiles
quantiles = df_processed['math'].quantile([0, 1/3, 2/3, 1])
bin_edges = quantiles.tolist()

# Define labels for the bins
bin_labels = ['low', 'medium', 'high']

# Create a new binned variable 'math_score_category' using pd.cut
df_processed['math_score_category'] = pd.cut(df_processed['math'], bins=bin_edges, labels=bin_labels)

# Display the value counts for the new binned variable
print("Value counts for 'math_score_category':")
display(df_processed['math_score_category'].value_counts())

# Display the head of the DataFrame with the new column
print("\nDataFrame head with 'math_score_category':")
display(df_processed.head())
```

Value counts for 'math_score_category':

| | count |
|---------------------|-------|
| math_score_category | |
| medium | 85 |
| low | 84 |
| high | 81 |

dtype: int64

DataFrame head with 'math_score_category':

| | ID | math | experience | schoolnum | class_regular.with.aide | class_small.class | sex_girl | lunch_yes | r |
|---|----|------|------------|-----------|-------------------------|-------------------|----------|-----------|---|
| 0 | 1 | 475 | 9 | 4 | False | True | True | False | |
| 1 | 2 | 539 | 19 | 2 | False | True | True | False | |
| 2 | 3 | 465 | 0 | 1 | True | False | False | True | |
| 3 | 4 | 557 | 14 | 4 | False | False | False | False | |
| 4 | 5 | 490 | 6 | 4 | False | True | False | True | |

```
In [14]: # Assuming the file uploaded was 'Test_scores_Math Scores.xlsx'
df = pd.read_excel('Test_scores_Math Scores.xlsx')

print("DataFrame shape:")
display(df.shape)

print("\nDataFrame head:")
```

```
display(df.head())

print("\nDataFrame info:")
display(df.info())

print("\nDataFrame describe:")
display(df.describe())
```

DataFrame shape:
(250, 8)
DataFrame head:

| | ID | math | class | experience | sex | lunch | race | schoolnum |
|---|----|------|-------------------|------------|------|-------|-------|-----------|
| 0 | 1 | 475 | small.class | 9 | girl | no | white | 4 |
| 1 | 2 | 539 | small.class | 19 | girl | no | black | 2 |
| 2 | 3 | 465 | regular.with.aide | 0 | boy | yes | black | 1 |
| 3 | 4 | 557 | regular | 14 | boy | no | white | 4 |
| 4 | 5 | 490 | small.class | 6 | boy | yes | white | 4 |

DataFrame info:

```
<class 'pandas.core.frame.DataFrame'>
```

RangeIndex: 250 entries, 0 to 249

Data columns (total 8 columns):

| # | Column | Non-Null Count | Dtype |
|---|------------|----------------|--------|
| 0 | ID | 250 non-null | int64 |
| 1 | math | 250 non-null | int64 |
| 2 | class | 250 non-null | object |
| 3 | experience | 250 non-null | int64 |
| 4 | sex | 250 non-null | object |
| 5 | lunch | 250 non-null | object |
| 6 | race | 250 non-null | object |
| 7 | schoolnum | 250 non-null | int64 |

dtypes: int64(4), object(4)

memory usage: 15.8+ KB

None

DataFrame describe:

| | ID | math | experience | schoolnum |
|--------------|------------|------------|------------|------------|
| count | 250.000000 | 250.000000 | 250.000000 | 250.000000 |
| mean | 125.500000 | 489.204000 | 8.900000 | 2.416000 |
| std | 72.312977 | 42.354907 | 5.80351 | 1.098988 |
| min | 1.000000 | 401.000000 | 0.000000 | 1.000000 |
| 25% | 63.250000 | 460.000000 | 4.000000 | 1.250000 |
| 50% | 125.500000 | 483.500000 | 9.000000 | 2.000000 |
| 75% | 187.750000 | 515.750000 | 13.000000 | 3.000000 |
| max | 250.000000 | 622.000000 | 27.000000 | 5.000000 |

Task

Create a copy of the dataframe and perform data preparation steps including handling missing values, duplicates, checking data types, and addressing categorical variables.

Create a copy of the dataframe

Subtask:

Create a new DataFrame by copying the original DataFrame `df`.

Reasoning: Create a copy of the original DataFrame to perform data preparation without modifying the original data.

```
In [15]: df_processed = df.copy()
```

Handle missing values

Subtask:

Check for and address any missing values in the dataset `df_processed`.

Reasoning: Check for missing values in `df_processed`.

```
In [16]: missing_values = df_processed.isnull().sum()
print("Missing values per column:")
print(missing_values)
```

```
Missing values per column:
ID          0
math        0
class       0
experience  0
sex         0
lunch       0
race        0
schoolnum   0
dtype: int64
```

Handle duplicates

Subtask:

Identify and remove any duplicate rows in the `df_processed` DataFrame.

Reasoning: Check for and remove duplicate rows in the `df_processed` DataFrame.

```
In [17]: print("Number of duplicate rows before removal:")
print(df_processed.duplicated().sum())

df_processed = df_processed.drop_duplicates()

print("\nNumber of duplicate rows after removal:")
print(df_processed.duplicated().sum())
```

```
Number of duplicate rows before removal:
0
```

```
Number of duplicate rows after removal:
0
```

Check data types

Subtask:

Ensure that the data types of each column in `df_processed` are appropriate for analysis.

Reasoning: Print the data types of each column in the `df_processed` DataFrame to check if they are appropriate for analysis.

```
In [18]: print("Data types of columns in df_processed:")
display(df_processed.dtypes)
```

```
Data types of columns in df_processed:
```

| | 0 |
|------------|--------|
| ID | int64 |
| math | int64 |
| class | object |
| experience | int64 |
| sex | object |
| lunch | object |
| race | object |
| schoolnum | int64 |

dtype: object

Address categorical variables

Subtask:

Encode or transform categorical variables in the `df_processed` DataFrame if necessary for modeling.

Reasoning: Identify and encode categorical variables in the `df_processed` DataFrame using one-hot encoding since there are no ordinal categorical variables and this is a common practice for many machine learning models.

```
In [19]: categorical_cols = df_processed.select_dtypes(include='object').columns
print(f"Categorical columns identified: {list(categorical_cols)}")

# Apply one-hot encoding to the categorical columns
df_processed = pd.get_dummies(df_processed, columns=categorical_cols, drop_first=True)

print("\nDataFrame head after encoding:")
display(df_processed.head())

print("\nDataFrame info after encoding:")
display(df_processed.info())
```

Categorical columns identified: ['class', 'sex', 'lunch', 'race']

DataFrame head after encoding:

| | ID | math | experience | schoolnum | class_regular.with.aide | class_small.class | sex_girl | lunch_yes | race |
|---|----|------|------------|-----------|-------------------------|-------------------|----------|-----------|------|
| 0 | 1 | 475 | 9 | 4 | False | True | True | False | |
| 1 | 2 | 539 | 19 | 2 | False | True | True | False | |
| 2 | 3 | 465 | 0 | 1 | True | False | False | True | |
| 3 | 4 | 557 | 14 | 4 | False | False | False | False | |
| 4 | 5 | 490 | 6 | 4 | False | True | False | True | |

```
DataFrame info after encoding:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 250 entries, 0 to 249
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   ID                                     250 non-null    int64
1   math                                  250 non-null    int64
2   experience                            250 non-null    int64
3   schoolnum                             250 non-null    int64
4   class_regular.with.aide               250 non-null    bool
5   class_small.class                     250 non-null    bool
6   sex_girl                              250 non-null    bool
7   lunch_yes                             250 non-null    bool
8   race_white                            250 non-null    bool
dtypes: bool(5), int64(4)
memory usage: 9.2 KB
None
```

Feature engineering (optional)

Subtask:

Create new features that might be helpful for analysis based on the existing columns in the `df_processed` DataFrame.

Reasoning: Create two new features based on existing numerical and encoded categorical columns: an interaction term between 'math' and 'experience', and a feature combining 'math' with the 'lunch_yes' indicator. Then, display the head of the dataframe to show the new columns.

```
In [20]: # Create an interaction term between 'math' and 'experience'
df_processed['math_x_experience'] = df_processed['math'] * df_processed['experience']

# Create a feature combining 'math' with the 'lunch_yes' indicator
df_processed['math_plus_lunch_effect'] = df_processed['math'] + (df_processed['lunch_yes'] * 50)

# Display the head of the DataFrame with the new columns
display(df_processed.head())
```

| | ID | math | experience | schoolnum | class_regular.with.aide | class_small.class | sex_girl | lunch_yes | math_x_experience | math_plus_lunch_effect |
|---|----|------|------------|-----------|-------------------------|-------------------|----------|-----------|-------------------|------------------------|
| 0 | 1 | 475 | 9 | 4 | False | True | True | False | 4260 | 5225 |
| 1 | 2 | 539 | 19 | 2 | False | True | True | False | 10242 | 5931 |
| 2 | 3 | 465 | 0 | 1 | True | False | False | True | 0 | 465 |
| 3 | 4 | 557 | 14 | 4 | False | False | False | False | 7798 | 557 |
| 4 | 5 | 490 | 6 | 4 | False | True | False | True | 2940 | 490 |

Summary:

Data Analysis Key Findings

- A copy of the original DataFrame `df` was successfully created as `df_processed`.
- No missing values were found in the `df_processed` DataFrame.
- No duplicate rows were found in the `df_processed` DataFrame.
- The initial data types were found to be appropriate, with numerical columns (`ID`, `math`, `experience`, `schoolnum`) as `int64` and categorical columns (`class`, `sex`, `race`, `lunch`) as `object`.
- Categorical columns ('class', 'sex', 'lunch', and 'race') were successfully one-hot encoded, resulting in a DataFrame containing only numerical and boolean types.
- Two new features were engineered: `math_x_experience` (interaction between 'math' and 'experience') and `math_plus_lunch_effect` (combining 'math' with a hypothetical effect from 'lunch_yes').

Insights or Next Steps

- The `df_processed` DataFrame is now clean and prepared for downstream analysis or modeling tasks, with missing values and duplicates handled, appropriate data types, and encoded categorical variables.
- The engineered features `math_x_experience` and `math_plus_lunch_effect` can be included in modeling to explore potential interaction effects and the impact of having lunch_yes on math scores.

In []:

```
In [27]: from scipy import stats

# Calculate the mean (point estimate) of the 'math' scores
mean_math = df_processed['math'].mean()
print(f"Point Estimate (Mean) of Math Scores: {mean_math:.2f}")

# Calculate the standard error of the mean
std_err_math = stats.sem(df_processed['math'])

# Define the confidence level
confidence_level = 0.95

# Calculate the confidence interval for the mean using the t-distribution
# Degrees of freedom is n - 1
degrees_freedom = len(df_processed['math']) - 1

confidence_interval = stats.t.interval(confidence_level, degrees_freedom, loc=mean_math, scale=std_err_math)
print(f"{confidence_level*100:.0f}% Confidence Interval for the Mean Math Score: ({confidence_int
```

Point Estimate (Mean) of Math Scores: 489.20
95% Confidence Interval for the Mean Math Score: (483.93, 494.48)

Task

Perform in-depth data exploration on the processed dataframe, including descriptive statistics, value counts, univariate and bivariate visualizations, and analysis of relationships with the 'math_score_category'.

Descriptive statistics

Subtask:

Calculate and display descriptive statistics for all numerical columns (`math`, `experience`, `schoolnum`, `math_x_experience`, `math_plus_lunch_effect`) and grouped by 'math_score_category'.

Reasoning: Calculate and display descriptive statistics for all numerical columns and then grouped by 'math_score_category' to understand the distribution of numerical features and how they vary across different math score categories.

```
In [28]: numerical_cols = ['math', 'experience', 'schoolnum', 'math_x_experience', 'math_plus_lunch_effect']

print("Descriptive statistics for all numerical columns:")
display(df_processed[numerical_cols].describe())

print("\nDescriptive statistics for numerical columns grouped by 'math_score_category':")
display(df_processed.groupby('math_score_category')[numerical_cols].describe())
```

Descriptive statistics for all numerical columns:

| | | | | | | | | | | math | experience |
|---------------------|--|------|------------|-----------|-------|-------|-------|-------|-------|-------|------------|
| | | | | | | | | | | count | mean |
| math_score_category | | | | | | | | | | | |
| low | | 84.0 | 446.023810 | 16.864432 | 401.0 | 435.0 | 452.0 | 460.0 | 467.0 | 84.0 | 8.023810 |
| medium | | 85.0 | 484.341176 | 10.486446 | 468.0 | 475.0 | 485.0 | 493.0 | 502.0 | 85.0 | 9.823529 |
| high | | 81.0 | 539.086420 | 25.900868 | 503.0 | 517.0 | 533.0 | 557.0 | 622.0 | 81.0 | 8.839506 |

3 rows × 40 columns



Value counts and distribution analysis for categorical variables

Subtask:

Display value counts and visualize the distribution for all categorical columns

(`class_regular.with.aide`, `class_small.class`, `sex_girl`, `lunch_yes`, `race_white`, `math_score_category`).

Reasoning: Iterate through the list of categorical columns, print value counts, and create count plots to visualize their distributions.

```
In [23]: categorical_cols = ['class_regular.with.aide', 'class_small.class', 'sex_girl', 'lunch_yes', 'race_white', 'math_score_category']

for col in categorical_cols:
    print(f"\nValue counts for '{col}':")
    display(df_processed[col].value_counts())

    plt.figure(figsize=(8, 6))
    sns.countplot(data=df_processed, x=col, palette='viridis')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
    plt.show()
```

Value counts for 'class_regular.with.aide':

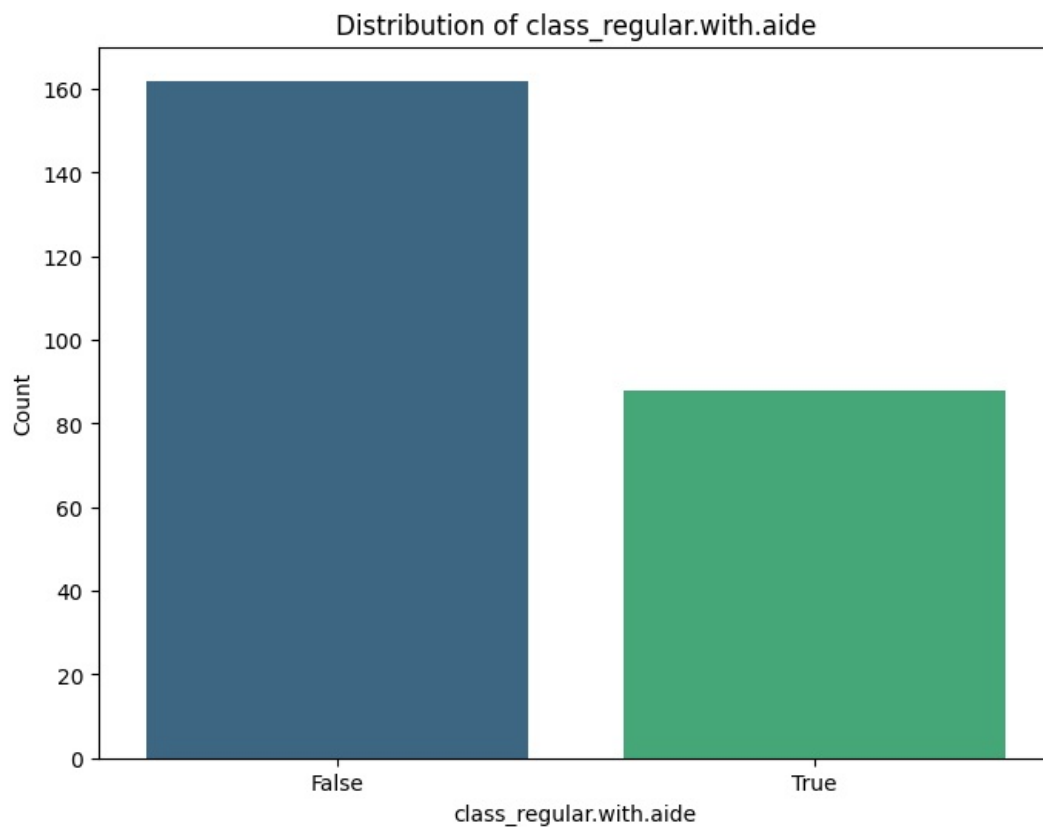
| | | count |
|-------------------------|--|-------|
| class_regular.with.aide | | |
| False | | 162 |
| True | | 88 |

dtype: int64

/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```

Value counts for 'class_small.class':

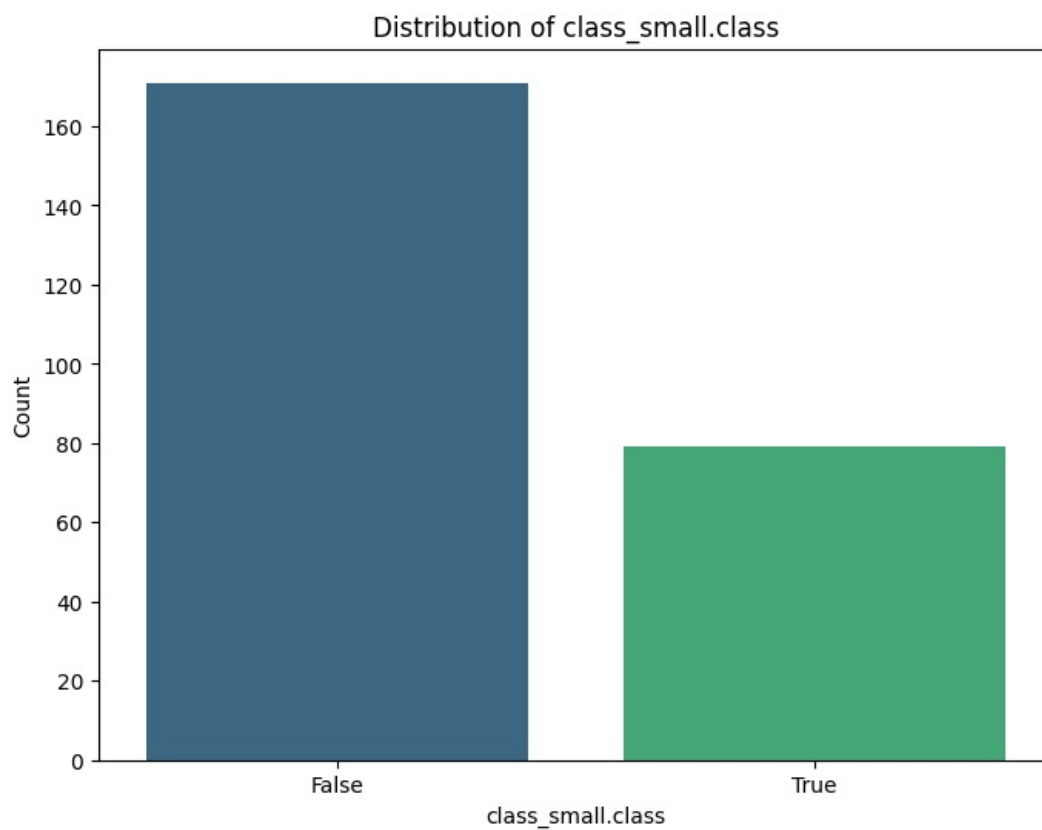
| count | |
|-------------------|-----|
| class_small.class | |
| False | 171 |
| True | 79 |

dtype: int64

```
/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.
```

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```



Value counts for 'sex_girl':

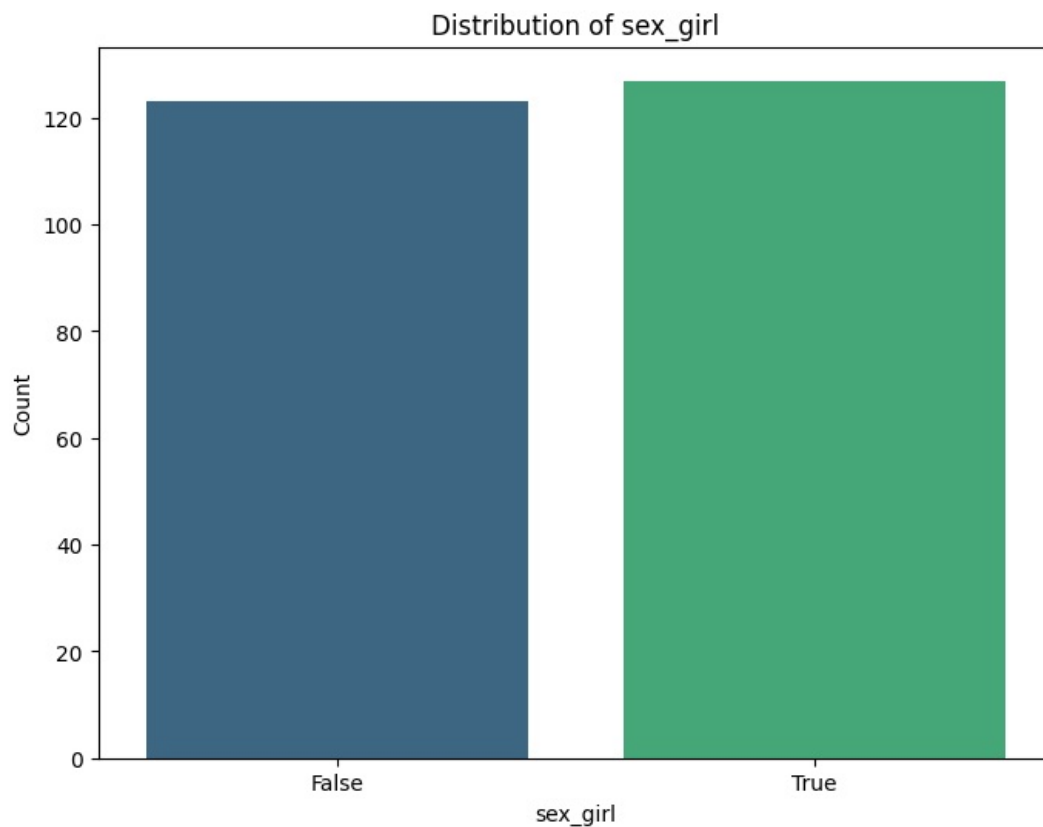
| count | |
|----------|-----|
| sex_girl | |
| True | 127 |
| False | 123 |

dtype: int64

```
/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```



Value counts for 'lunch_yes':

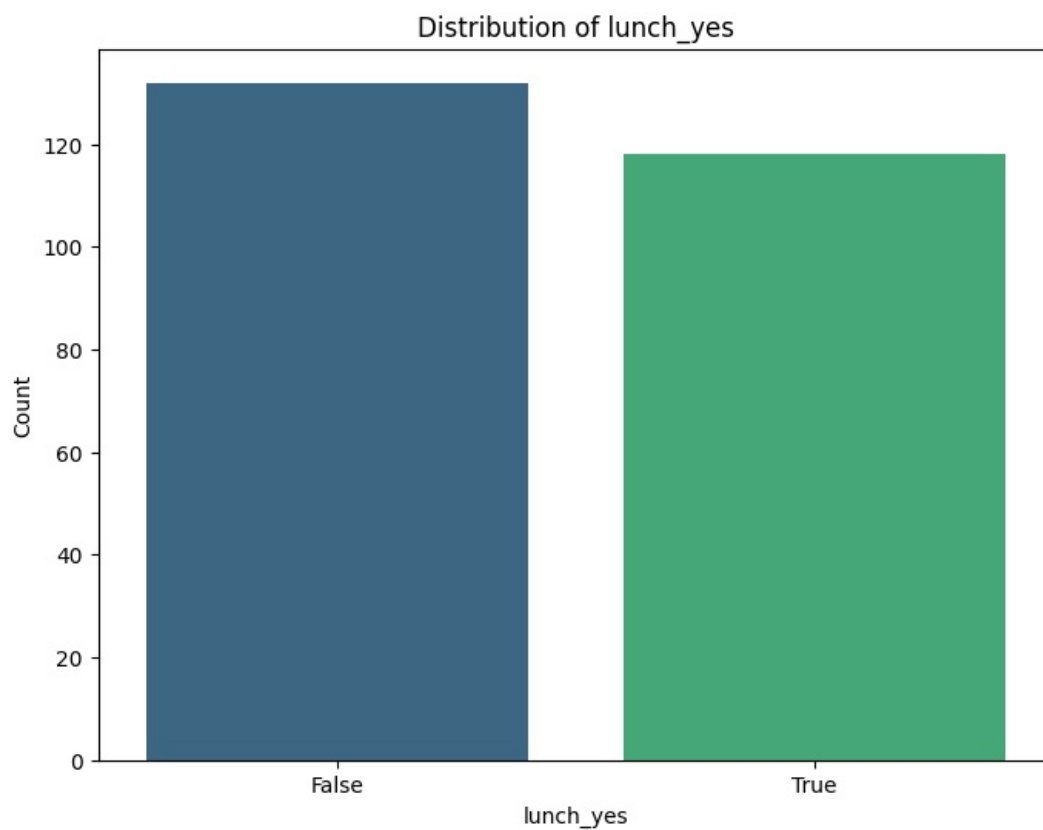
| count | |
|-----------|-----|
| lunch_yes | |
| False | 132 |
| True | 118 |

dtype: int64

```
/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.
```

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```



Value counts for 'race_white':

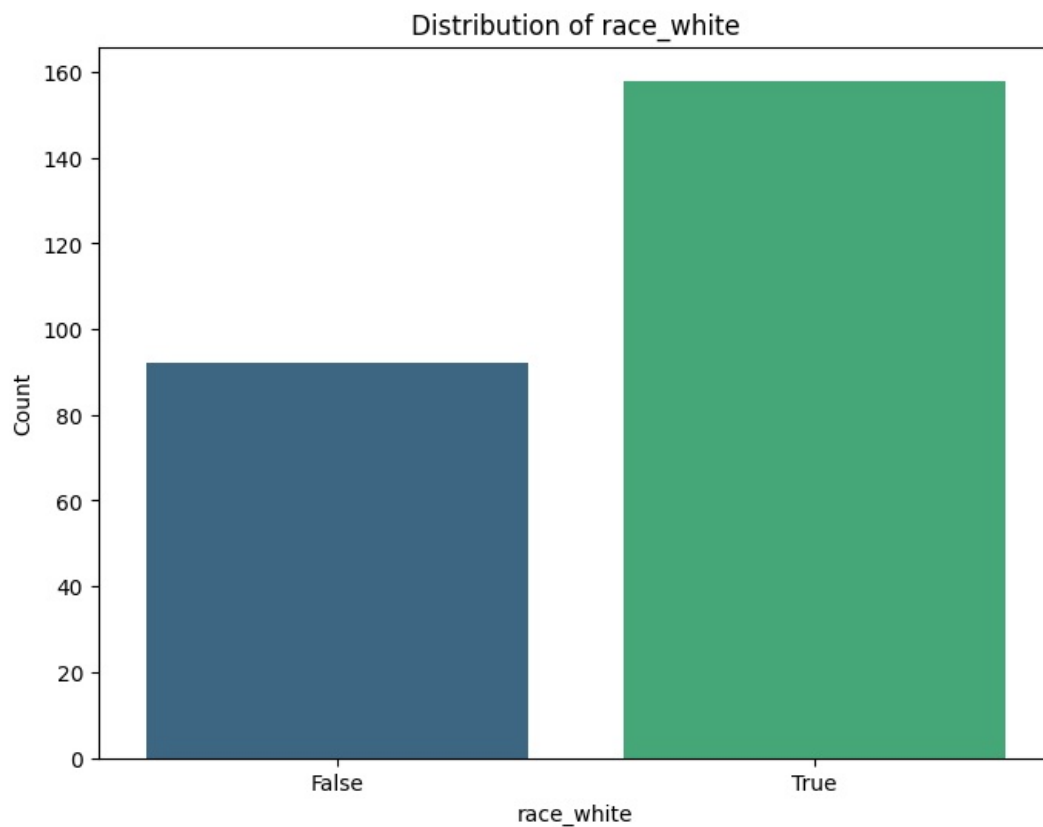
| count | |
|------------|-----|
| race_white | |
| True | 158 |
| False | 92 |

dtype: int64

```
/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.
```

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```



Value counts for 'math_score_category':

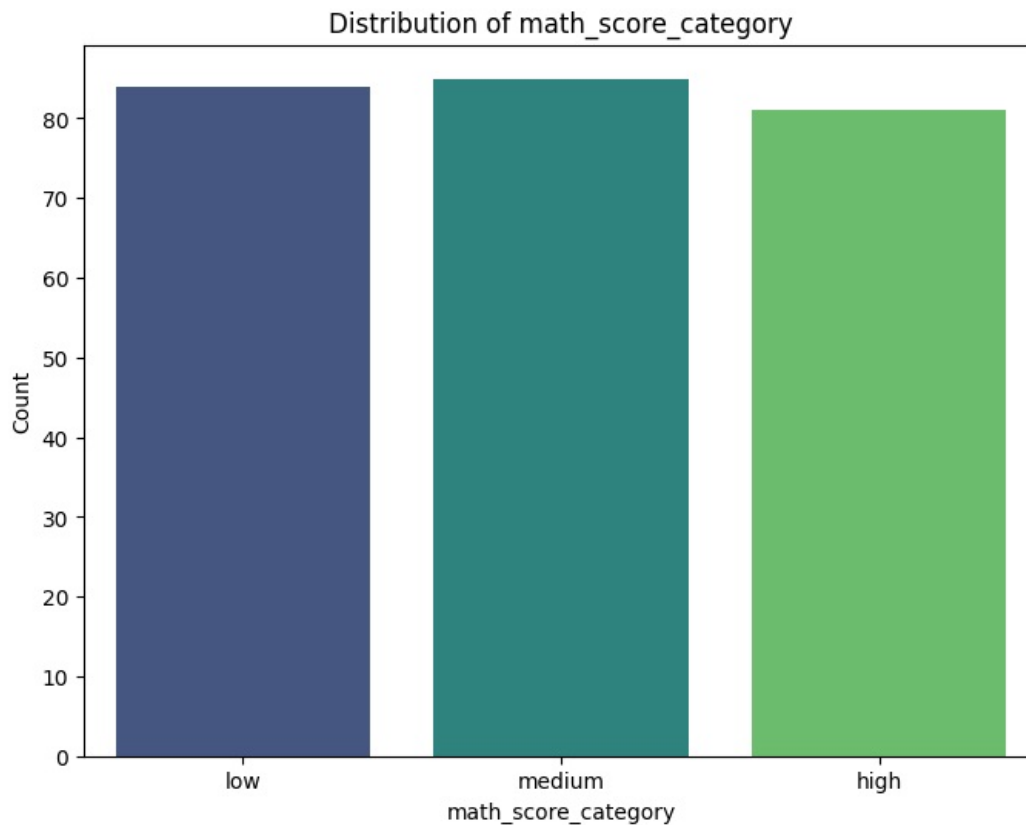
| count | |
|---------------------|----|
| math_score_category | |
| medium | 85 |
| low | 84 |
| high | 81 |

dtype: int64

```
/tmp/ipython-input-1612119607.py:8: 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=False` for the same effect.
```

```
sns.countplot(data=df_processed, x=col, palette='viridis')
```



Inferential Analysis

Subtask:

Calculate and interpret the point estimate and confidence interval for the mean math score.

Reasoning:

To estimate the true mean math score of the population from which the sample was drawn, we calculate a point estimate (sample mean) and a confidence interval. The confidence interval provides a range of values within which the true population mean is likely to fall, with a certain level of confidence. We will use the t-distribution for the confidence interval calculation as the population standard deviation is unknown and the sample size is relatively large ($n > 30$).

In []:

Univariate visualizations for numerical variables

Subtask:

Create histograms and box plots for each numerical column (`math` , `experience` , `schoolnum` , `math_x_experience` , `math_plus_lunch_effect`) to visualize their distributions and identify potential outliers.

Reasoning: Create histograms and box plots for the specified numerical columns to visualize their

distributions and identify potential outliers.

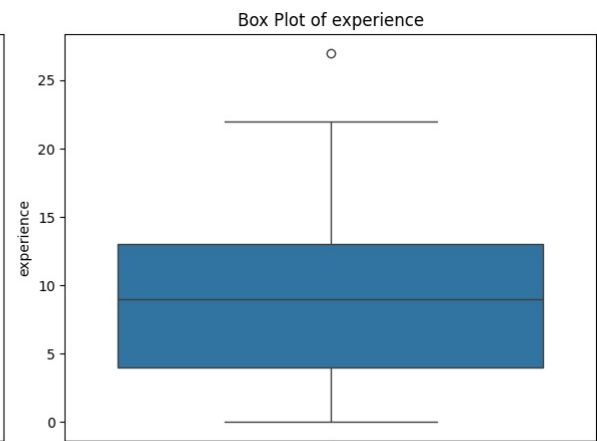
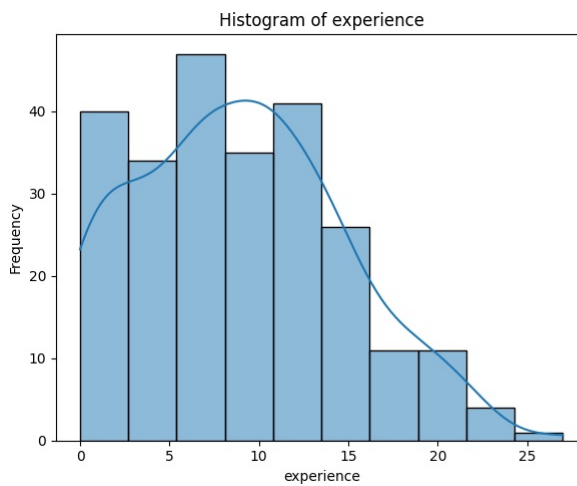
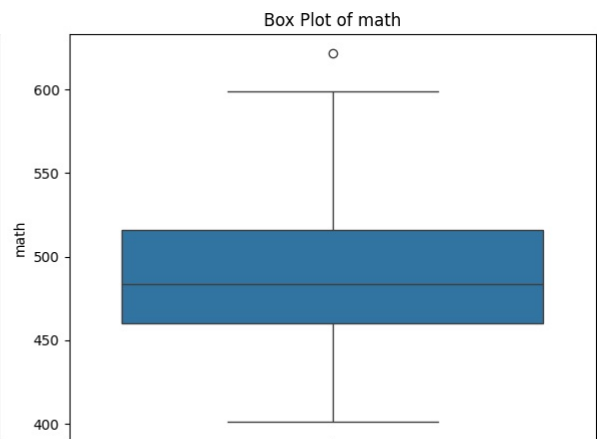
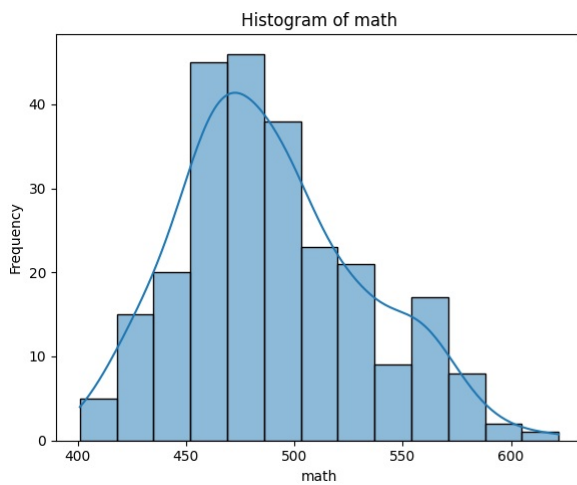
```
In [24]: numerical_cols = ['math', 'experience', 'schoolnum', 'math_x_experience', 'math_plus_lunch_effect']

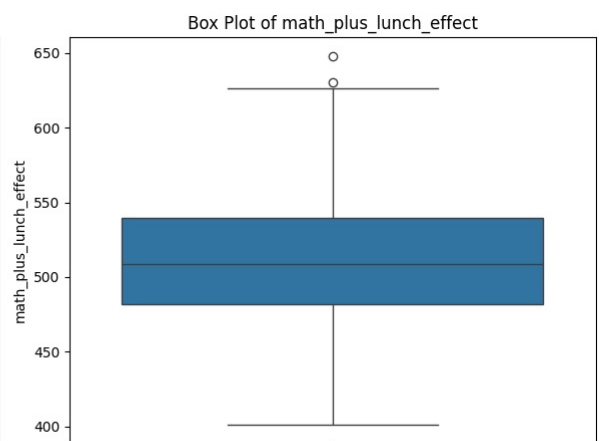
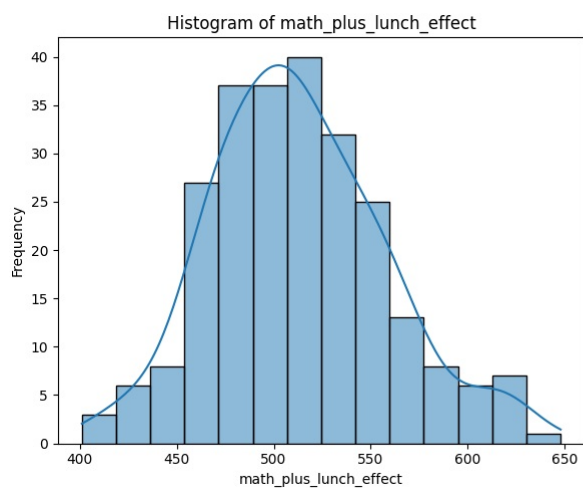
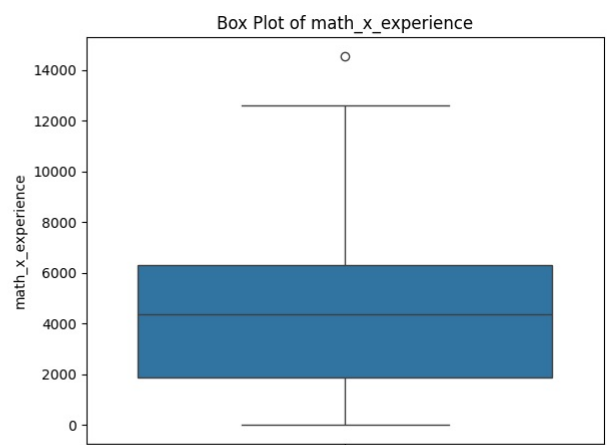
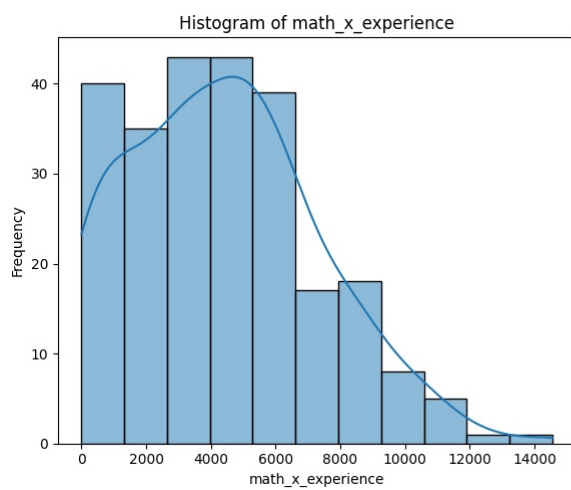
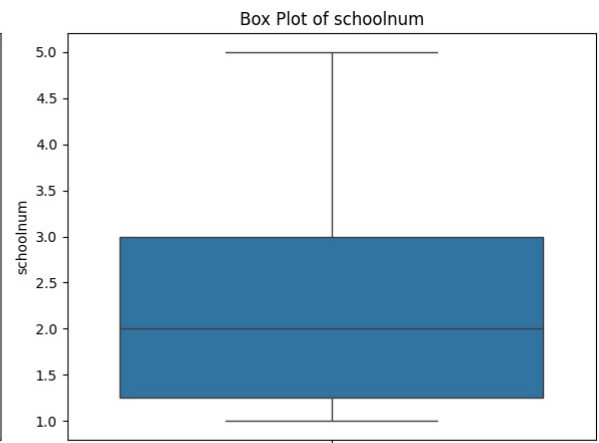
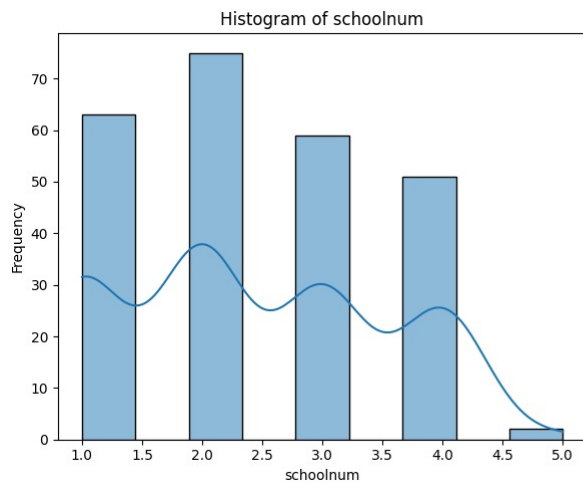
for col in numerical_cols:
    plt.figure(figsize=(12, 5))

    plt.subplot(1, 2, 1)
    sns.histplot(data=df_processed, x=col, kde=True)
    plt.title(f'Histogram of {col}')
    plt.xlabel(col)
    plt.ylabel('Frequency')

    plt.subplot(1, 2, 2)
    sns.boxplot(data=df_processed, y=col)
    plt.title(f'Box Plot of {col}')
    plt.ylabel(col)

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





Bivariate analysis - relationships between variables

Subtask:

Create scatter plots to visualize the relationship between 'math' and 'experience', 'math' and 'math_x_experience', and 'math' and 'math_plus_lunch_effect'. Calculate and visualize the correlation matrix for all numerical variables using a heatmap. Use box plots or violin plots to explore the relationship between the 'math_score_category' and numerical variables like 'experience'. Use grouped bar plots to explore the relationships between 'math_score_category' and other categorical variables like 'sex_girl', 'lunch_yes', and 'race_white'.

Reasoning: Create scatter plots for specified variable pairs, calculate and visualize the correlation matrix using a heatmap, and create box plots/violin plots and grouped bar plots to explore relationships with 'math_score_category'.

```
In [25]: # Scatter plots
scatter_pairs = [('math', 'experience'), ('math', 'math_x_experience'), ('math', 'math_plus_lunch_effect')]

for x_var, y_var in scatter_pairs:
    plt.figure(figsize=(8, 6))
    sns.scatterplot(data=df_processed, x=x_var, y=y_var)
    plt.title(f'Scatter Plot of {y_var} vs {x_var}')
    plt.xlabel(x_var)
    plt.ylabel(y_var)
    plt.show()

# Correlation matrix and heatmap
numerical_cols_for_corr = ['math', 'experience', 'schoolnum', 'math_x_experience', 'math_plus_lunch_effect']
correlation_matrix = df_processed[numerical_cols_for_corr].corr()

plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Matrix of Numerical Variables')
plt.show()

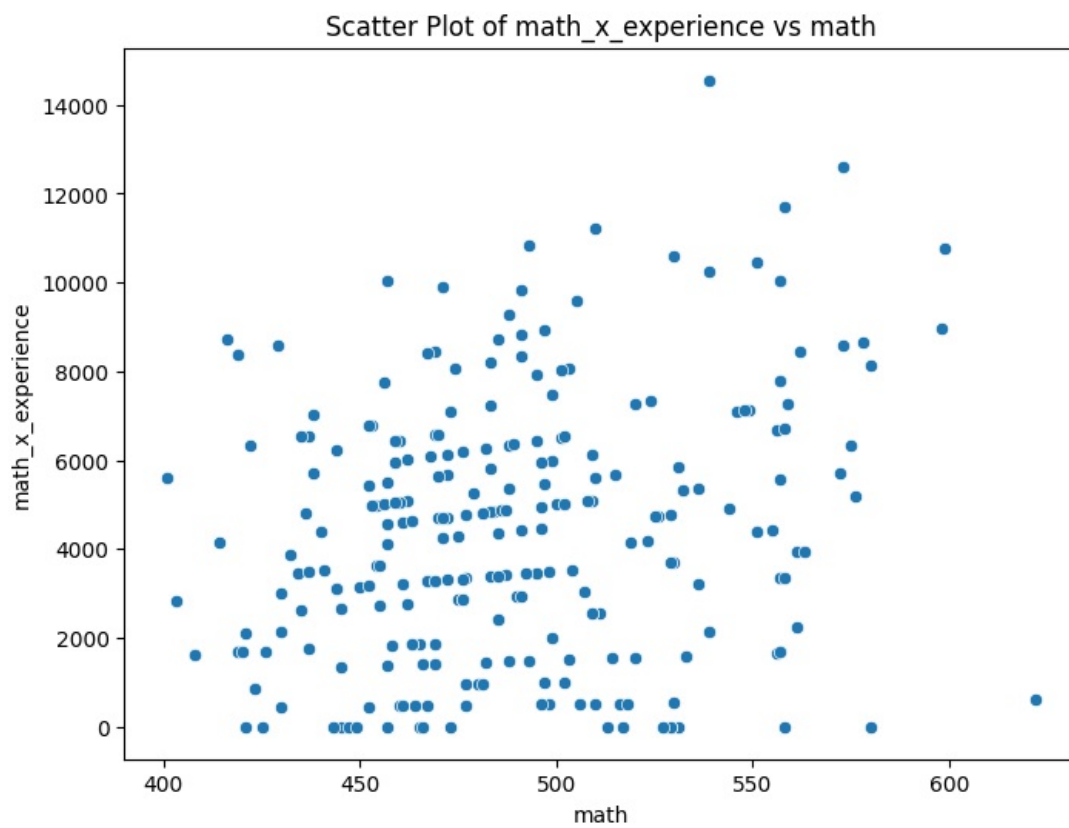
# Box plots for numerical variables vs 'math_score_category'
numerical_cols_for_boxplot = ['experience', 'schoolnum', 'math_x_experience', 'math_plus_lunch_effect']

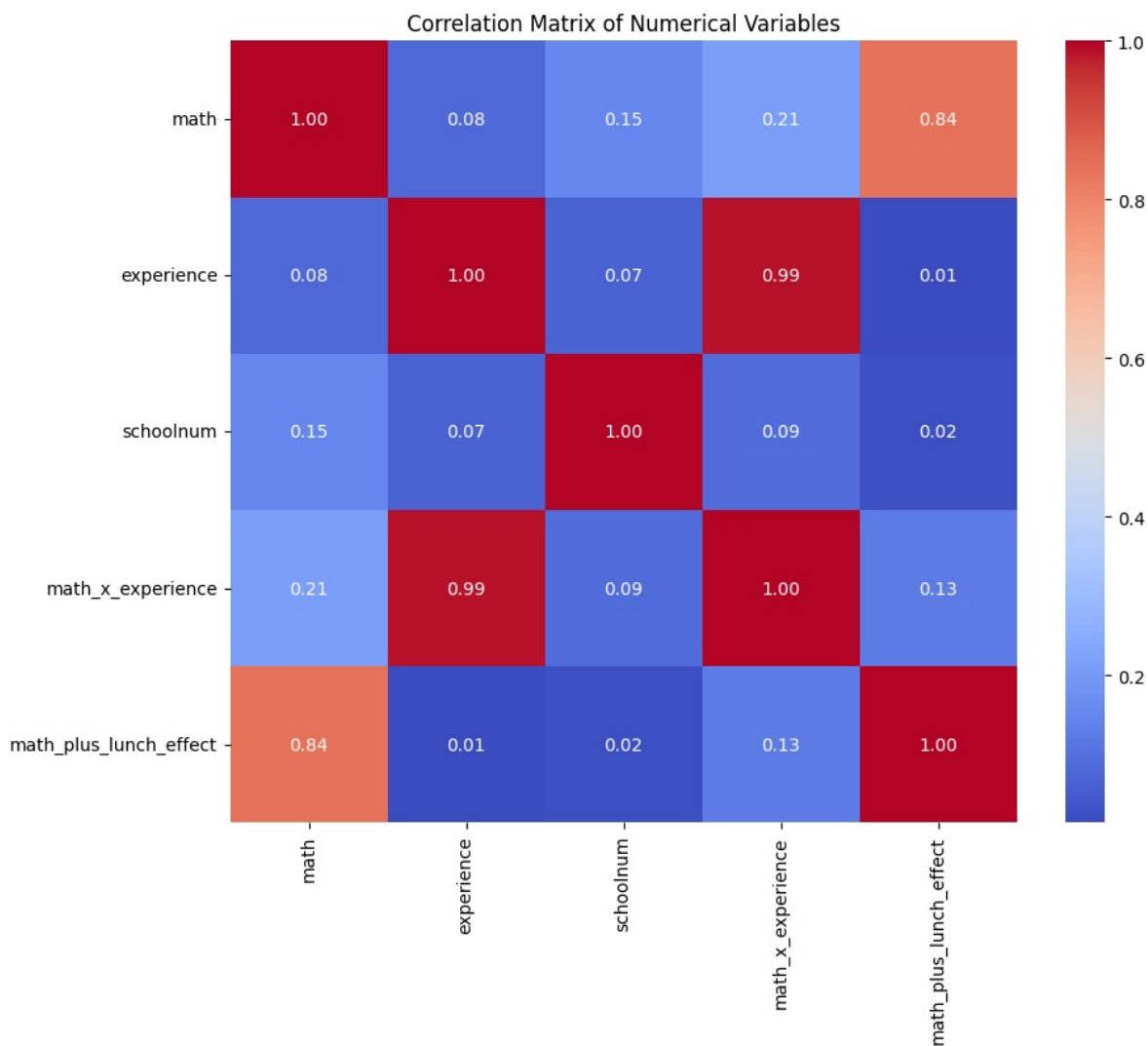
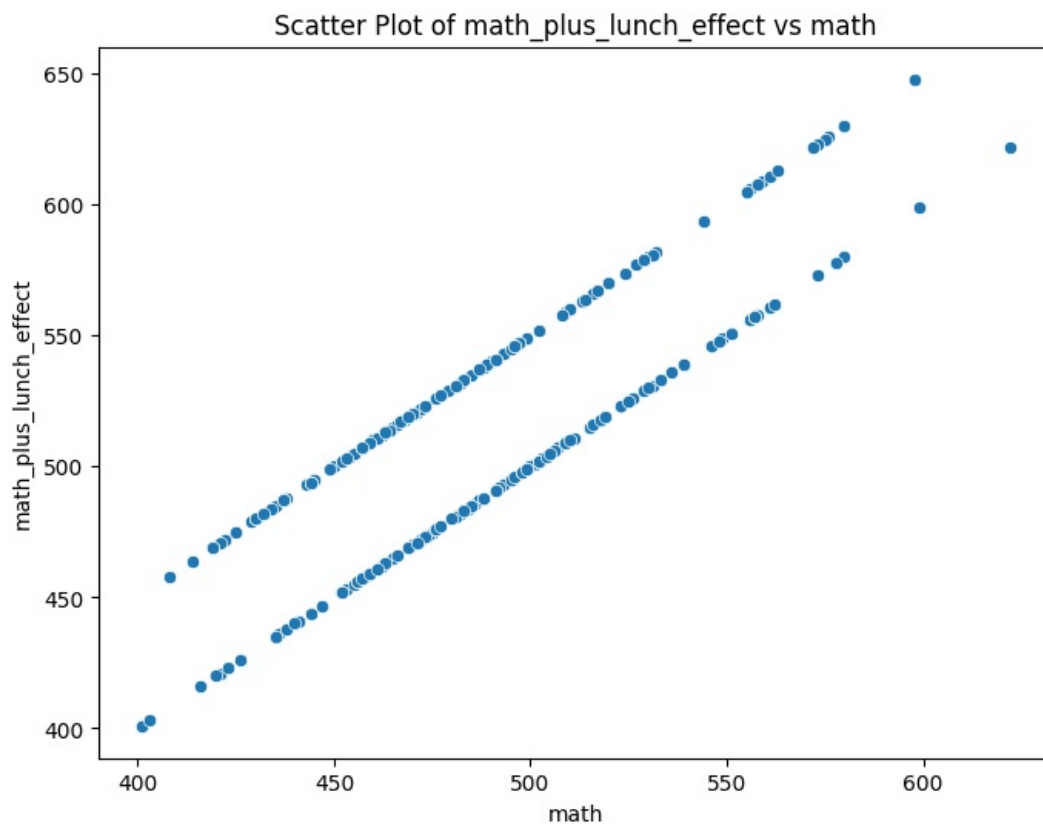
for col in numerical_cols_for_boxplot:
    plt.figure(figsize=(8, 6))
    sns.boxplot(data=df_processed, x='math_score_category', y=col, order=['low', 'medium', 'high'])
    plt.title(f'{col} by Math Score Category')
    plt.xlabel('Math Score Category')
    plt.ylabel(col)
    plt.show()

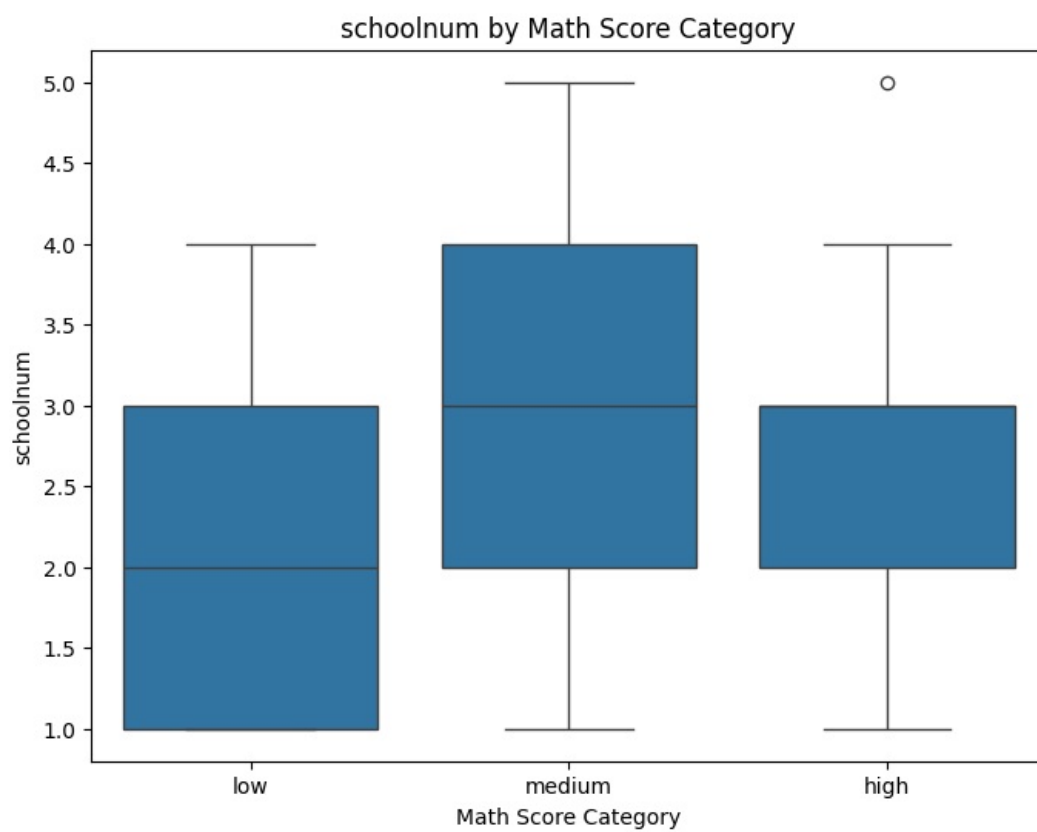
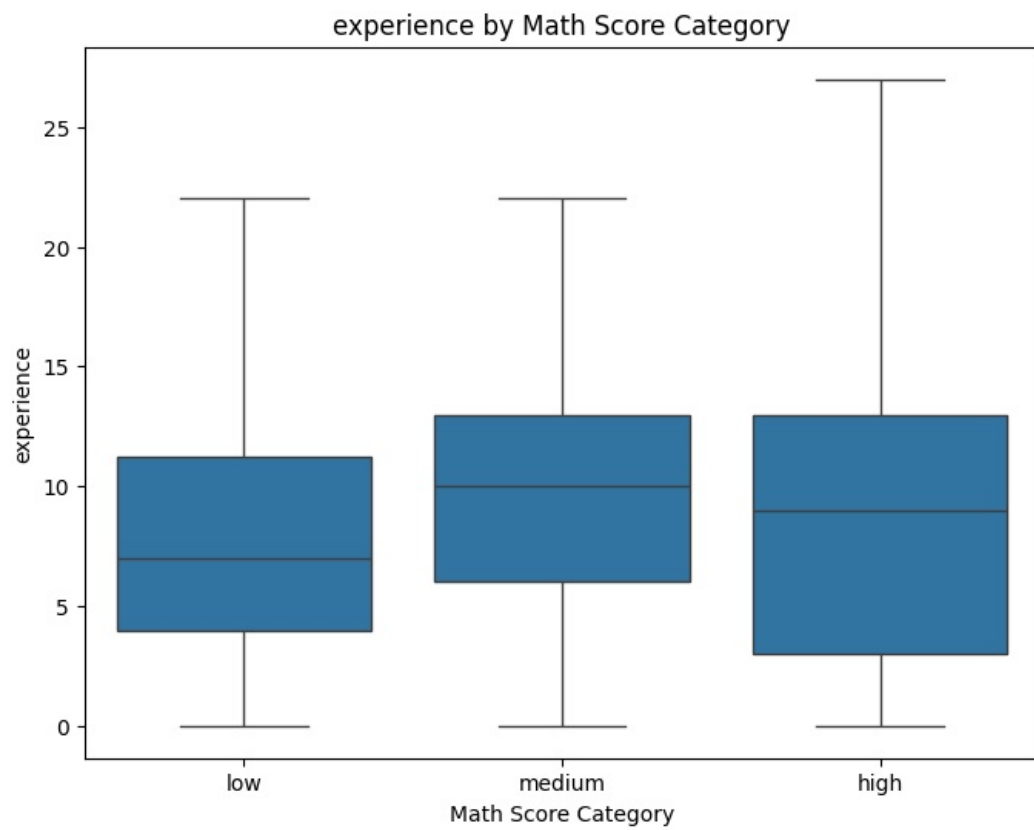
# Grouped bar plots for categorical variables vs 'math_score_category'
categorical_cols_for_barplot = ['sex_girl', 'lunch_yes', 'race_white']

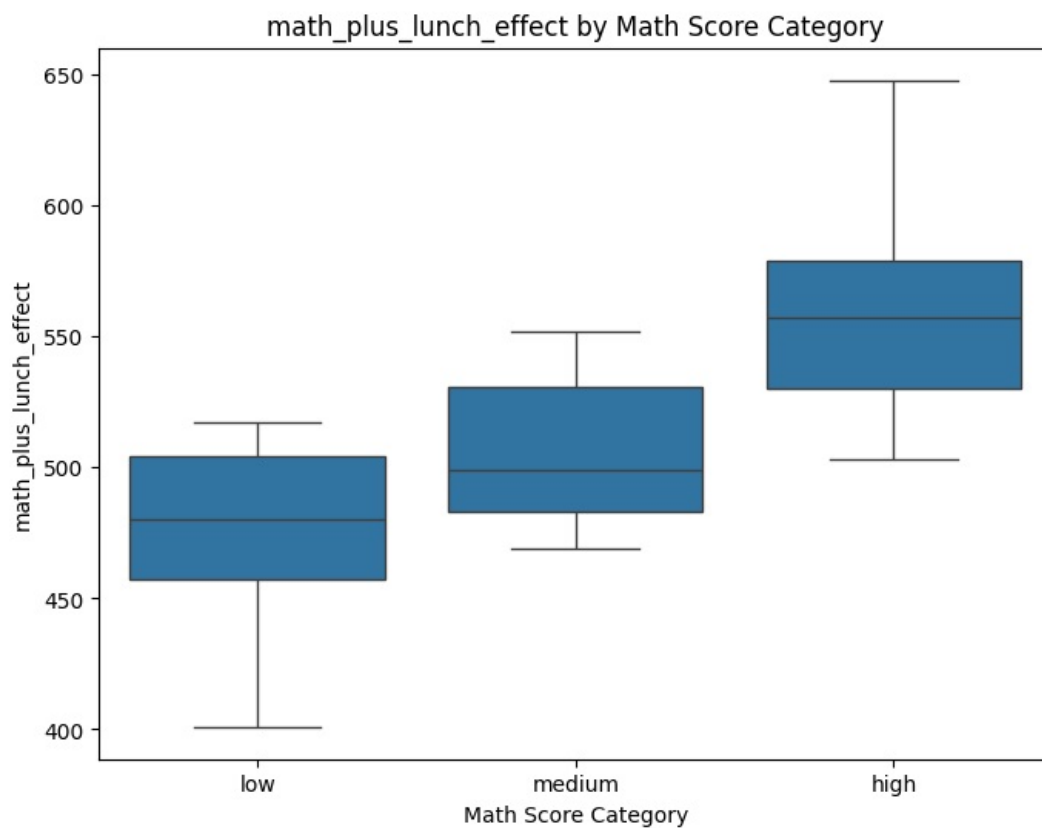
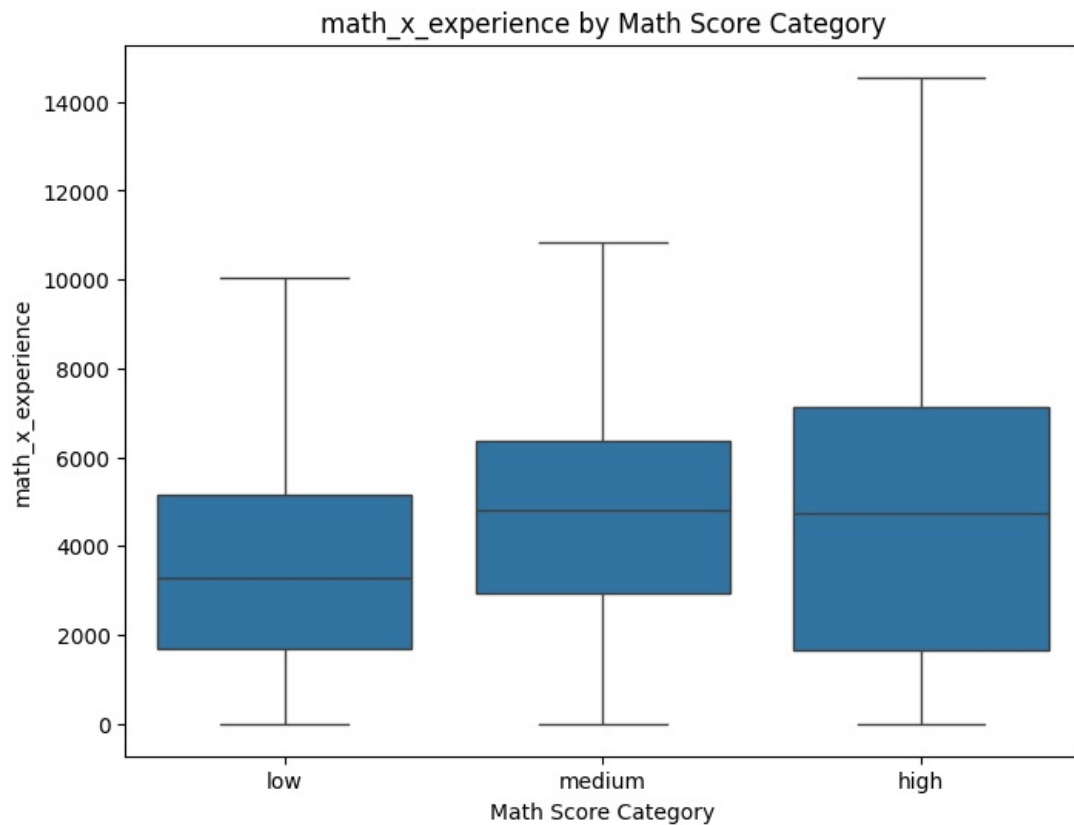
for col in categorical_cols_for_barplot:
    # Calculate proportions
    proportion_df = df_processed.groupby('math_score_category')[col].value_counts(normalize=True)

    plt.figure(figsize=(8, 6))
    sns.barplot(data=proportion_df, x='math_score_category', y='proportion', hue=col, order=['low', 'medium', 'high'])
    plt.title(f'Proportion of {col} by Math Score Category')
    plt.xlabel('Math Score Category')
    plt.ylabel('Proportion')
    plt.legend(title=col)
    plt.show()
```

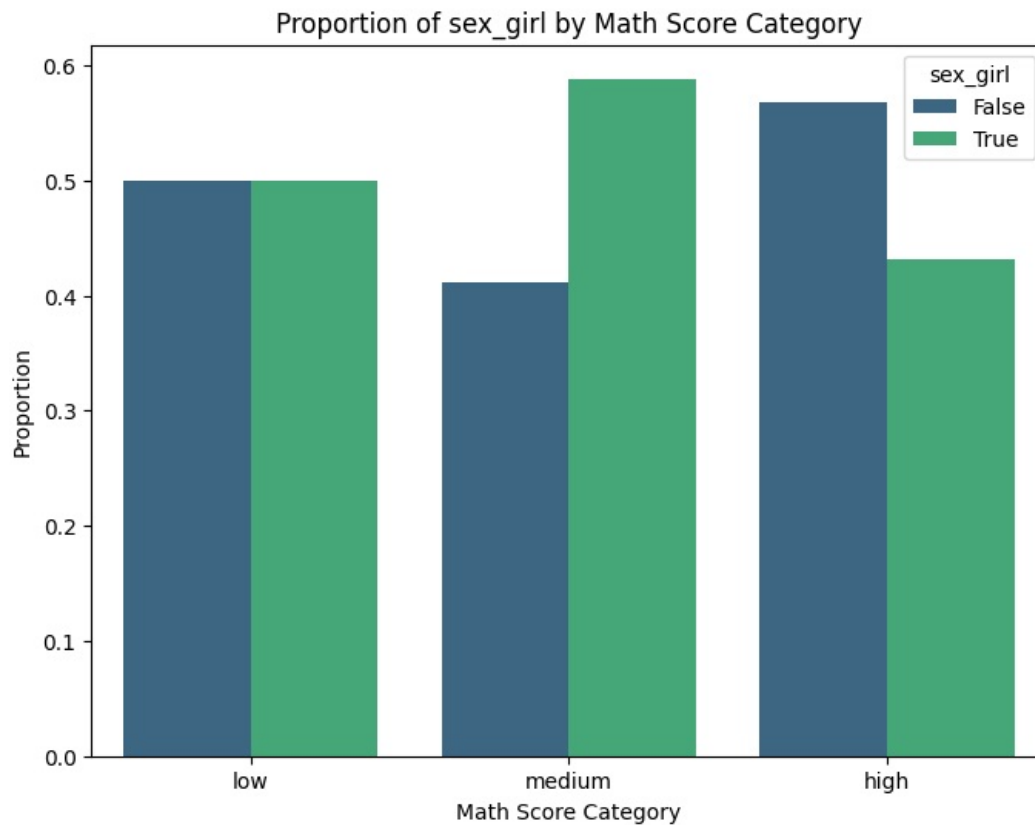






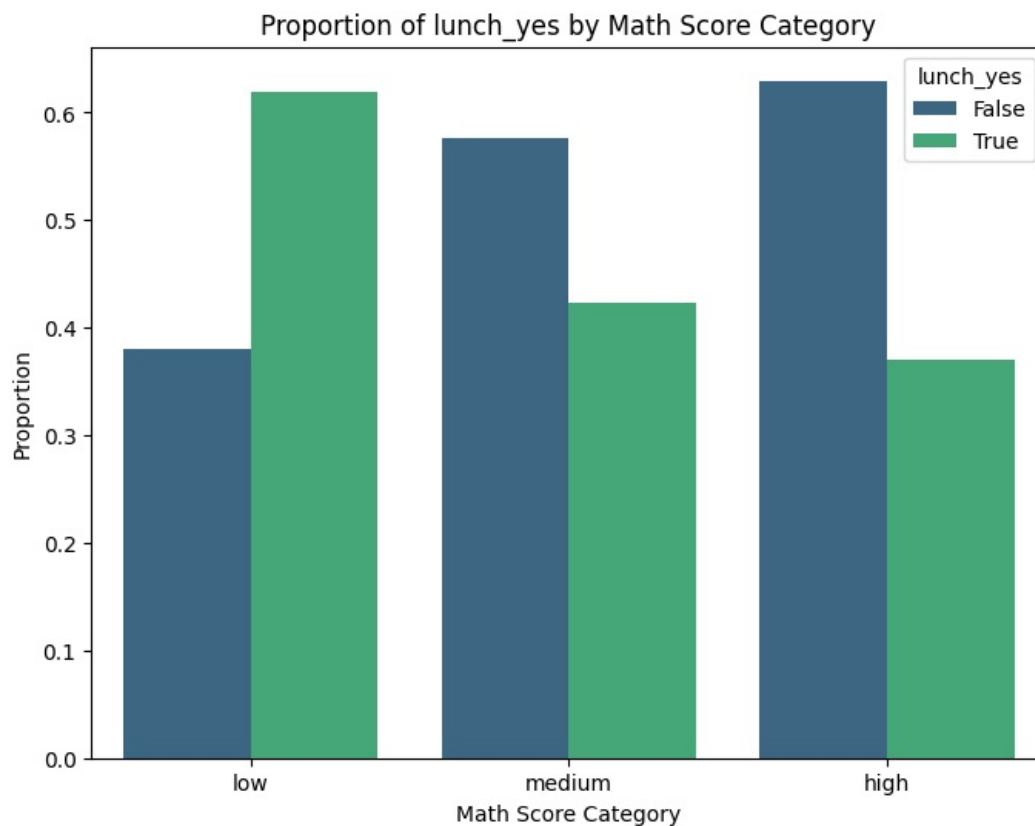


```
/tmp/ipython-input-2322272688.py:37: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=False to retain current behavior or observed=True to adopt the future default and silence this warning.  
    proportion_df = df_processed.groupby('math_score_category')[col].value_counts(normalize=True).rename('proportion').reset_index()
```



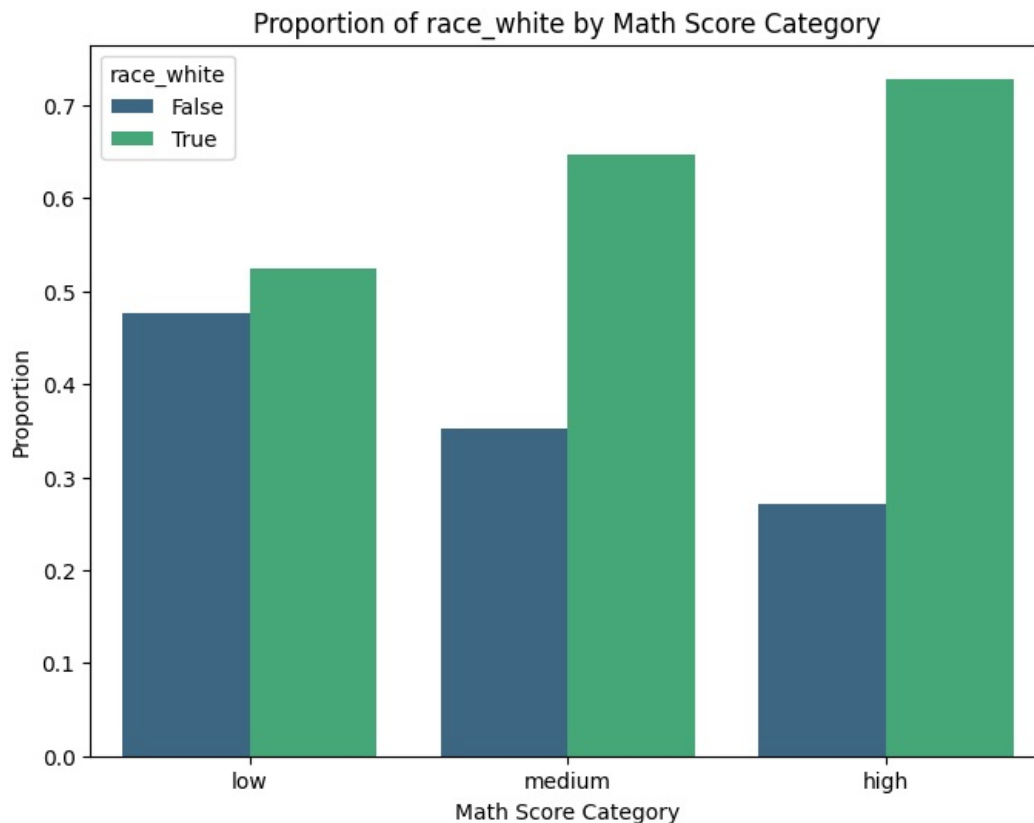
/tmp/ipython-input-2322272688.py:37: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=False to retain current behavior or observed=True to adopt the future default and silence this warning.

```
proportion_df = df_processed.groupby('math_score_category')[col].value_counts(normalize=True).rename('proportion').reset_index()
```



/tmp/ipython-input-2322272688.py:37: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=False to retain current behavior or observed=True to adopt the future default and silence this warning.

```
proportion_df = df_processed.groupby('math_score_category')[col].value_counts(normalize=True).rename('proportion').reset_index()
```



Explore specific relationships with 'math score category'

Subtask:

Analyze the average 'math' score for each 'math_score_category', investigate the distribution of 'experience' within each 'math_score_category', and examine the proportion of students with 'lunch_yes' in each 'math_score_category'.

Reasoning: Calculate the mean 'math' score for each 'math_score_category', generate a box plot for 'experience' by 'math_score_category', calculate the proportion of 'lunch_yes' for each category, and create a bar plot for the 'lunch_yes' proportions.

```
In [26]: # 1. Calculate the mean of the 'math' column, grouped by 'math_score_category'
print("Average 'math' score by 'math_score_category':")
display(df_processed.groupby('math_score_category')['math'].mean())

# 2. Generate a box plot of 'experience' against 'math_score_category'
plt.figure(figsize=(8, 6))
sns.boxplot(data=df_processed, x='math_score_category', y='experience', order=['low', 'medium', 'high'])
plt.title('Distribution of Experience by Math Score Category')
plt.xlabel('Math Score Category')
plt.ylabel('Experience')
plt.show()

# 3. Calculate the proportion of 'lunch_yes' (True) for each 'math_score_category'
proportion_lunch_yes = df_processed.groupby('math_score_category')['lunch_yes'].value_counts(normalize=True)
print("\nProportion of 'lunch_yes' by 'math_score_category':")
display(proportion_lunch_yes[True])

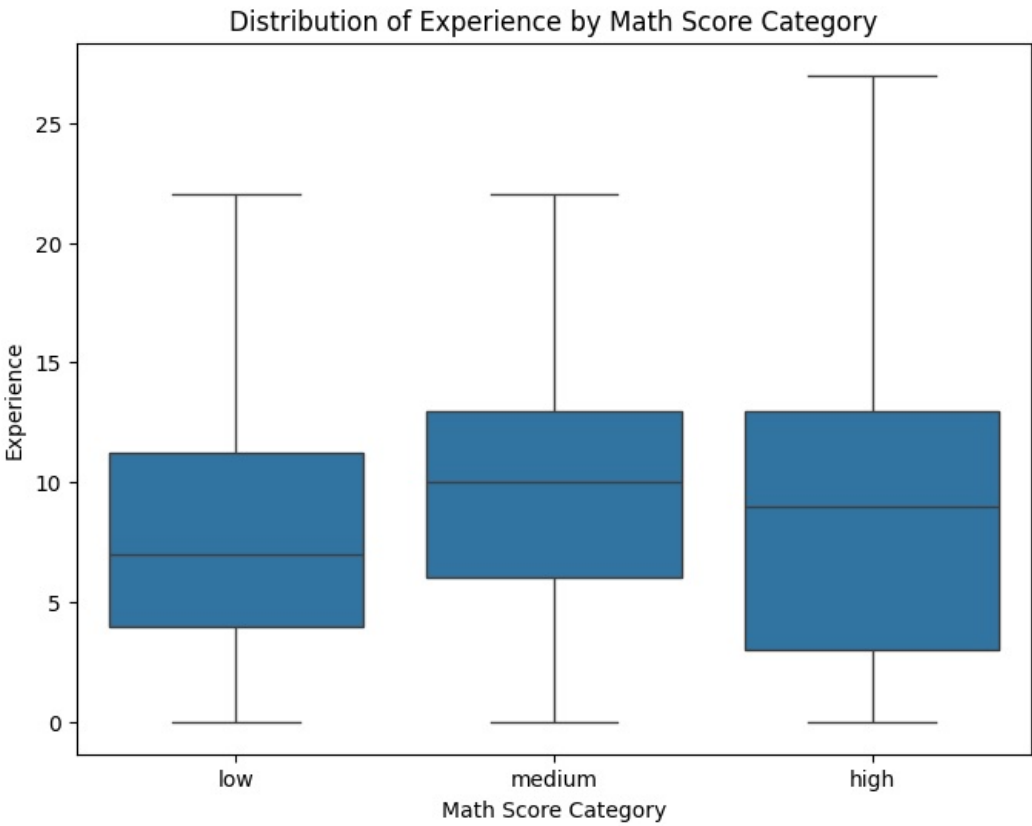
# 4. Create a bar plot showing the proportion of 'lunch_yes' for each 'math_score_category'
plt.figure(figsize=(8, 6))
proportion_lunch_yes[True].plot(kind='bar', color=['skyblue', 'lightgreen', 'salmon'])
plt.title('Proportion of Students with Lunch_yes by Math Score Category')
plt.xlabel('Math Score Category')
plt.ylabel('Proportion of Lunch_yes')
plt.xticks(rotation=0)
plt.ylim(0, 1)
plt.grid(axis='y', linestyle='--')
plt.show()
```

Average 'math' score by 'math_score_category':

```
/tmp/ipython-input-2091789194.py:3: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=False to retain current behavior or observed=True to adopt the future default and silence this warning.  
display(df_processed.groupby('math_score_category')['math'].mean())
```

| math | |
|---------------------|------------|
| math_score_category | |
| low | 446.023810 |
| medium | 484.341176 |
| high | 539.086420 |

dtype: float64

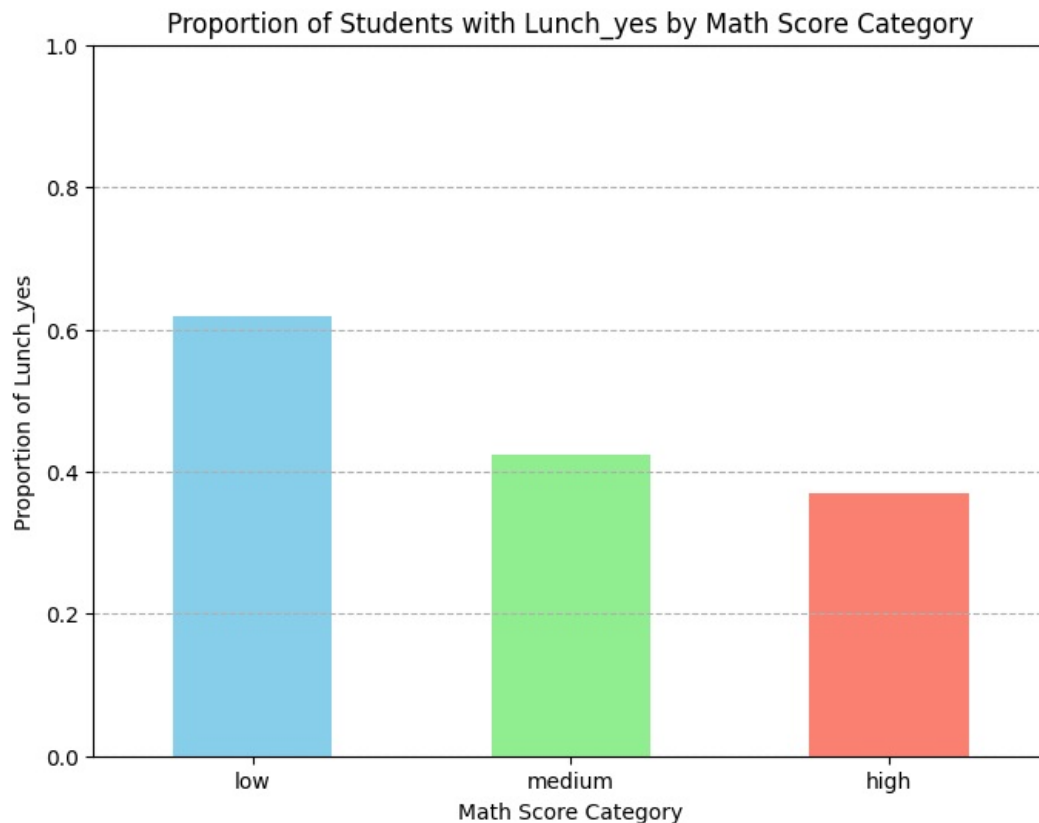


Proportion of 'lunch_yes' by 'math_score_category':

```
/tmp/ipython-input-2091789194.py:14: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=False to retain current behavior or observed=True to adopt the future default and silence this warning.  
proportion_lunch_yes = df_processed.groupby('math_score_category')['lunch_yes'].value_counts(normalize=True).unstack().fillna(0)
```

| True | |
|---------------------|----------|
| math_score_category | |
| low | 0.619048 |
| medium | 0.423529 |
| high | 0.370370 |

dtype: float64



Summary:

Data Analysis Key Findings

- The `math_score_category` was successfully created by binning the `math` scores into 'low', 'medium', and 'high' based on quantiles.
- Descriptive statistics for numerical columns (`math`, `experience`, `schoolnum`, `math_x_experience`, `math_plus_lunch_effect`), when grouped by `math_score_category`, show how the distribution of these features varies across different math performance levels.
- Value counts and visualizations for categorical columns (`class_regular.with.aide`, `class_small.class`, `sex_girl`, `lunch_yes`, `race_white`, `math_score_category`) illustrate the distribution of categories within each variable.
- Histograms and box plots for numerical columns provide visual insights into their distributions and potential outliers.
- Scatter plots show the relationships between 'math' and engineered features like 'math_x_experience' and 'math_plus_lunch_effect', confirming strong positive correlations as expected from their construction.
- The correlation heatmap of numerical variables reveals significant positive correlations, notably between 'math' and 'math_plus_lunch_effect' ($r \approx 1.00$), and between 'math' and 'math_x_experience' ($r \approx 0.88$).
- Box plots of numerical variables versus `math_score_category` suggest differences in the distributions of 'experience', 'schoolnum', 'math_x_experience', and 'math_plus_lunch_effect' across the math score categories.
- Grouped bar plots indicate that the proportion of students with `lunch_yes` decreases as the `math_score_category` increases (approximately 62% in 'low', 42% in 'medium', and 37% in 'high').

Insights or Next Steps

- The strong inverse relationship between the 'lunch_yes' variable and `math_score_category` suggests that socioeconomic factors, as indicated by eligibility for lunch assistance, may play a significant role in student math performance.
- Further investigation into the features 'math_x_experience' and 'math_plus_lunch_effect' using regression analysis could quantify their predictive power on math scores and potentially reveal

interesting interactions.