```
import pandas as pd
import math
import plotly.express as px
from scipy.stats import levene, bartlett, f_oneway, norm
import scikit_posthocs as sp
```

- **Imports:** The script begins by importing necessary libraries:

    o  pandas for data manipulation.

    o  math for mathematical operations.

    o  plotly.express for data visualization (not used in this part of the script).

    o  scipy.stats for statistical tests (Levene's test, Bartlett's test, ANOVA, and normal distribution functions).

    o  scikit_posthocs for conducting Dunn's post-hoc test.

```
# Load the dataset
file_path = 'C:\\Users\\loydt\\Downloads\\Projects\\Superstore Sales
Dataset.xlsx'
data = pd.read_excel(file_path)
```

- **Loading Data:** The script loads the Superstore Sales Dataset from a specified file path into a DataFrame called data.

```
# Filter the big four states generating higher sales
states_of_interest = ['Washington', 'California', 'New York', 'Florida',
'Pennsylvania']
state_data = data[data['State'].isin(states_of_interest)]
```

- **Filtering Data:** It filters the data to include only entries from specific states known for higher sales. The filtered data is stored in state_data.

```
# Extract relevant columns and create a deep copy to avoid
SettingWithCopyWarning
high_sales_states = state_data[['Segment', 'State', 'City', 'Region', 'Ship
Mode', 'Order Date', 'Category', 'Sub-Category', 'Product Name',
'Sales']].copy()
```

- **Data Preparation:** Relevant columns related to sales analysis are extracted, and a deep copy is created to avoid the SettingWithCopyWarning when modifying the DataFrame later.

```
# Take logarithm of the Sales column, ensuring non-positive values are
handled
high_sales_states['log_sales'] = high_sales_states['Sales'].apply(lambda x:
math.log(x) if x > 0 else None)
```

- **Log Transformation:** A logarithmic transformation is applied to the `Sales` column to reduce skewness and stabilize variance, ensuring non-positive values are handled appropriately.

```python
# Check for null values in 'log_sales' after transformation
if high_sales_states['log_sales'].isnull().any():
    print("Warning: There are non-positive sales values that have been transformed to NaN.")
```

- **Null Check:** The script checks for any null values in the `log_sales` column resulting from the transformation. If found, it prints a warning message.

```python
# Create a DataFrame with 'log_sales' and 'Category'
dunn_data = high_sales_states[['log_sales', 'Category']].dropna()  # Drop NaN values for valid analysis
```

- **Prepare Data for Dunn's Test:** A new DataFrame, `dunn_data`, is created with only the `log_sales` and `Category` columns, dropping any rows with NaN values to ensure valid analysis.

```python
# Conduct Dunn's Test
dunn_results = sp.posthoc_dunn(dunn_data, val_col='log_sales', group_col='Category', p_adjust='bonferroni')
```

- **Dunn's Test Execution:** Dunn's post-hoc test is performed using the `posthoc_dunn` function from `scikit_posthocs`. This test is suitable for pairwise comparisons following a Kruskal-Wallis test. The p-values are adjusted using the Bonferroni method to control for type I error.

```python
# Convert the results to a DataFrame
dunn_results_df = dunn_results.stack().reset_index()
dunn_results_df.columns = ['group1', 'group2', 'p-adj']
```

- **Result Formatting:** The results from Dunn's test are converted into a DataFrame, `dunn_results_df`, with appropriately named columns for group comparisons and adjusted p-values.

```python
# Calculate mean difference, lower and upper confidence intervals, and reject
dunn_results_df['meandiff'] = dunn_results_df.apply(
    lambda row: dunn_data[dunn_data['Category'] ==
row['group1']]['log_sales'].mean() -
                dunn_data[dunn_data['Category'] ==
row['group2']]['log_sales'].mean(), axis=1)
```

- **Mean Difference Calculation:** The script calculates the mean difference in log sales between each pair of categories and stores it in the `meandiff` column of the results DataFrame.

```python
# Create a 'reject' column based on the adjusted p-value
dunn_results_df['reject'] = dunn_results_df['p-adj'] < 0.05
```

- **Hypothesis Testing:** A new column, `reject`, is created to indicate whether to reject the null hypothesis for each comparison based on whether the adjusted p-value is less than 0.05.

```python
# Calculate confidence intervals
for index, row in dunn_results_df.iterrows():
    group1_data = dunn_data[dunn_data['Category'] ==
row['group1']]['log_sales']
    group2_data = dunn_data[dunn_data['Category'] ==
row['group2']]['log_sales']

    mean1 = group1_data.mean()
    mean2 = group2_data.mean()

    # Calculate standard error
    se1 = group1_data.std() / math.sqrt(len(group1_data))
    se2 = group2_data.std() / math.sqrt(len(group2_data))

    # Calculate mean difference
    meandiff = mean1 - mean2

    # Z-score for 95% confidence interval
    z_score = norm.ppf(0.975)  # 1.96 for 95% CI

    # Calculate confidence intervals
    lower_bound = meandiff - z_score * math.sqrt(se1**2 + se2**2)
    upper_bound = meandiff + z_score * math.sqrt(se1**2 + se2**2)

    # Assign lower and upper bounds to the DataFrame
    dunn_results_df.at[index, 'lower'] = lower_bound
    dunn_results_df.at[index, 'upper'] = upper_bound
```

- **Confidence Interval Calculation:** For each pair of categories, the script calculates the confidence intervals for the mean differences. It computes the mean and standard error for each group, uses the Z-score for a 95% confidence level, and then calculates the lower and upper bounds of the confidence interval. These values are stored in the `lower` and `upper` columns of the results DataFrame.

```python
# Display the results
print(dunn_results_df)
```

- **Output:** Finally, the script prints the DataFrame containing the results of Dunn's test, including the group comparisons, p-values, mean differences, rejection of null hypothesis, and confidence intervals.

This explanation summarizes the purpose and functionality of each part of your script, making it easier to understand the analysis performed on the Superstore Sales Dataset.