

Week 1 & 2: NumPy Operations and Tasks

Create NumPy arrays from Python Data Structures, Intrinsic NumPy objects and Random Functions.

```
# 1. Create NumPy arrays
# a. From Python data structures
import numpy as np

list_data = [1, 2, 3, 4, 5]
array_from_list = np.array(list_data)

# b. From intrinsic NumPy objects
zeros_array = np.zeros((3, 3))
ones_array = np.ones((2, 4))

# c. Using random functions
random_array = np.random.rand(3, 3)
print("Array from List:\n", array_from_list)
print("Zeros Array:\n", zeros_array)
print("Random Array:\n", random_array)

Array from List:
[1 2 3 4 5]
Zeros Array:
[[0. 0. 0.]
 [0. 0. 0.]
 [0. 0. 0.]]
Random Array:
[[0.66504087 0.4886642 0.72990495]
 [0.41879849 0.49981234 0.98057609]
 [0.90091381 0.90376282 0.12051952]]
```

Manipulation of NumPy arrays- Indexing, Slicing, Reshaping, Joining and Splitting.

```
# Indexing
print("First element:", array_from_list[0])

# Slicing
print("Slice of array:", array_from_list[1:4])

# Reshaping
reshaped_array = np.arange(1, 13).reshape(3, 4)
print("Reshaped Array:\n", reshaped_array)

# Joining
joined_array = np.concatenate((array_from_list, array_from_list))
```

```

print("Joined Array:\n", joined_array)

# Splitting
split_array = np.split(joined_array, 2)
print("Split Array:\n", split_array)

First element: 1
Slice of array: [2 3 4]
Reshaped Array:
[[ 1  2  3  4]
 [ 5  6  7  8]
 [ 9 10 11 12]]
Joined Array:
[1 2 3 4 5 1 2 3 4 5]
Split Array:
[array([1, 2, 3, 4, 5]), array([1, 2, 3, 4, 5])]

```

Computation on NumPy arrays using Universal Functions and Mathematical methods.

```

squared_array = np.square(array_from_list)
mean_value = np.mean(array_from_list)
print("Squared Array:\n", squared_array)
print("Mean Value:", mean_value)

Squared Array:
[ 1  4  9 16 25]
Mean Value: 3.0

```

Import a CSV file and perform various Statistical and Comparison operations on rows/columns.

```

import pandas as pd
data = {
'A': [10, 20, 30],
'B': [5, 15, 25],
'C': [8, 18, 28]
}
df = pd.DataFrame(data)
csv_path = "sample.csv"
df.to_csv(csv_path, index=False)

# Load the CSV file
csv_data = pd.read_csv(csv_path)
print("CSV Data:\n", csv_data)

# Statistical operations
column_sum = csv_data['A'].sum()
row_mean = csv_data.iloc[0].mean()
print("Sum of Column A:", column_sum)
print("Mean of Row 0:", row_mean)

```

```

# Comparison operations
greater_than_15 = csv_data > 15
print("Elements greater than 15:\n", greater_than_15)

CSV Data:
    A    B    C
0   105   8
1   20    15   18
2   30    25   28
Sum of Column A: 60
Mean of Row 0: 7.666666666666667
Elements greater than 15:
    A    B    C
0  False  False  False
1  True   False  True
2  True   True   True

```

Load an image file and do crop and flip operation using NumPy Indexing

```

from PIL import Image
import matplotlib.pyplot as plt
image_path = "C:\\\\Users\\\\B11202016\\\\Desktop\\\\nba.jpg" # Replace with
your image path
image = Image.open(image_path)
# Convert image to NumPy array
image_array = np.array(image)
# Crop operation
cropped_image = image_array[50:150, 50:150]
# Flip operation
flipped_image = np.flipud(image_array)

# Convert arrays back to images and save
cropped_image_pil = Image.fromarray(cropped_image)
flipped_image_pil = Image.fromarray(flipped_image)
cropped_image_pil.save("cropped_image.jpg")
plt.imshow(cropped_image_pil)
flipped_image_pil.save("flipped_image.jpg")
plt.imshow(flipped_image_pil)
print("Image operations completed. Cropped and flipped images saved.")

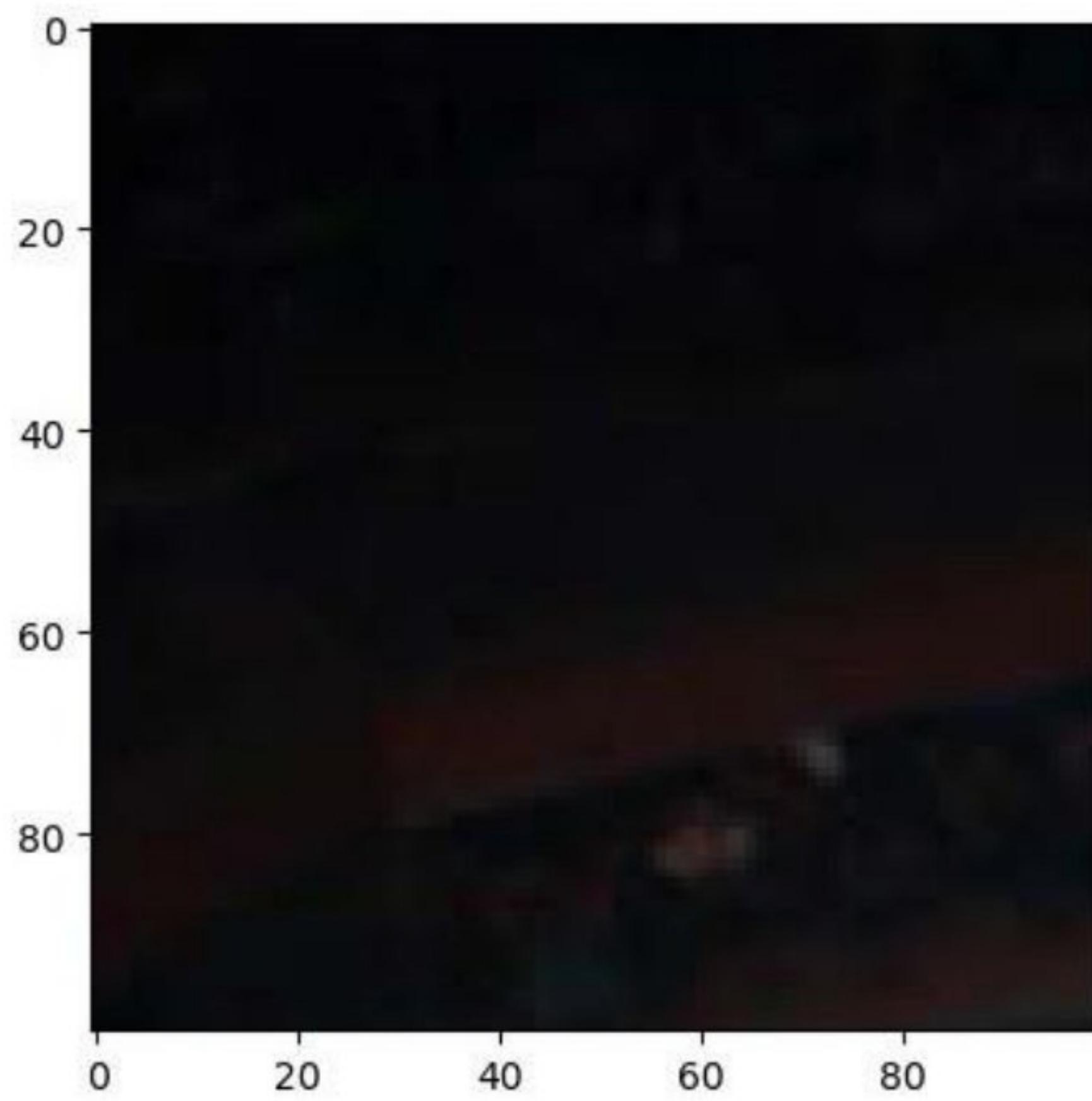
Image operations completed. Cropped and flipped images saved.

```



```
#image cropping
top,bottom,left,right=100,400,200,500
pic2=flipped_image_pil.crop((top,bottom,left,right))
plt.imshow(pic2)

<matplotlib.image.AxesImage at 0x20236a25670>
```



Week 3, 4, 5: Pandas Operations and Tasks

Create Pandas Series and Data Frame from various inputs.

```
# a. From a list
series_from_list = pd.Series([10, 20, 30, 40])
print("Series from List:\n", series_from_list)

# b. From a dictionary
df_from_dict = pd.DataFrame({
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [25, 30, 35],
    'Score': [85.5, 92.0, 78.0]
})
print("DataFrame from Dictionary:\n", df_from_dict)

Series from List:
0    10
1    20
2    30
3    40
dtype: int64
DataFrame from Dictionary:
   Name  Age  Score
```

```

0    Alice   25   85.5
1      Bob   30   92.0
2  Charlie   35   78.0

```

Import any CSV file to Pandas Data Frame and perform the following:

- (a) Visualize the first and last 10 records.
- (b) Get the shape, index and column details.
- (c) Select/Delete the records (rows)/columns based on conditions.
- (d) Perform ranking and sorting operations.
- (e) Do required statistical operations on the given columns.
- (f) Find the count and uniqueness of the given categorical values.
- (g) Rename single/multiple columns.

```

# Create a sample CSV
data = {
    'ID': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
    'Category': ['A', 'B', 'A', 'C', 'B', 'C', 'A', 'A', 'B', 'C'],
    'Value': [50, 60, 45, 70, 80, 75, 85, 90, 55, 65]
}
csv_file = "sample_week3.csv"
pd.DataFrame(data).to_csv(csv_file, index=False)

# Load the CSV
df = pd.read_csv(csv_file)

# a. Visualize the first and last 10 records
print("First 10 records:\n", df.head(10))
print("Last 10 records:\n", df.tail(10))

# b. Get the shape, index, and column details
print("Shape of DataFrame:", df.shape)
print("Index:", df.index)
print("Columns:", df.columns)

# c. Select/Delete records based on conditions

# Select rows where Value > 60
selected_rows = df[df['Value'] > 60]
print("Rows where Value > 60:\n", selected_rows)

# Delete rows where Category == 'B'
filtered_df = df[df['Category'] != 'B']
print("DataFrame after deleting rows where Category == 'B':\n",
      filtered_df)

```

```

# d. Perform ranking and sorting
df['Value_Rank'] = df['Value'].rank(ascending=False)
sorted_df = df.sort_values(by='Value', ascending=False)
print("Ranked DataFrame:\n", df)
print("Sorted DataFrame:\n", sorted_df)

# e. Statistical operations
mean_value = df['Value'].mean()
sum_value = df['Value'].sum()
print("Mean of Value column:", mean_value)
print("Sum of Value column:", sum_value)

# f. Count and uniqueness of categorical values
category_counts = df['Category'].value_counts()
unique_categories = df['Category'].unique()
print("Category Counts:\n", category_counts)
print("Unique Categories:\n", unique_categories)

# g. Rename single/multiple columns
df.rename(columns={'Value': 'Score_Value', 'Category': 'Group'},
inplace=True)
print("Renamed DataFrame:\n", df)

First 10 records:
   ID Category  Value
0   1         A     50
1   2         B     60
2   3         A     45
3   4         C     70
4   5         B     80
5   6         C     75
6   7         A     85
7   8         A     90
8   9         B     55
9  10        C     65

Last 10 records:
   ID Category  Value
0   1         A     50
1   2         B     60
2   3         A     45
3   4         C     70
4   5         B     80
5   6         C     75
6   7         A     85
7   8         A     90
8   9         B     55
9  10        C     65

Shape of DataFrame: (10, 3)
Index: RangeIndex(start=0, stop=10, step=1)

```

```

Columns: Index(['ID', 'Category', 'Value'], dtype='object')
Rows where Value > 60:
    ID Category  Value
3   4         C    70
4   5         B    80
5   6         C    75
6   7         A    85
7   8         A    90
9  10        C    65
DataFrame after deleting rows where Category == 'B':
    ID Category  Value
0   1         A    50
2   3         A    45
3   4         C    70
5   6         C    75
6   7         A    85
7   8         A    90
9  10        C    65
Ranked DataFrame:
    ID Category  Value  Value_Rank
0   1         A    50      9.0
1   2         B    60      7.0
2   3         A    45     10.0
3   4         C    70      5.0
4   5         B    80      3.0
5   6         C    75      4.0
6   7         A    85      2.0
7   8         A    90      1.0
8   9         B    55      8.0
9  10        C    65      6.0
Sorted DataFrame:
    ID Category  Value  Value_Rank
7   8         A    90      1.0
6   7         A    85      2.0
4   5         B    80      3.0
5   6         C    75      4.0
3   4         C    70      5.0
9  10        C    65      6.0
1   2         B    60      7.0
8   9         B    55      8.0
0   1         A    50      9.0
2   3         A    45     10.0
Mean of Value column: 67.5
Sum of Value column: 675
Category Counts:
  Category
A        4
B        3
C        3

```

```
Name: count, dtype: int64
```

```
Unique Categories:
```

```
['A' 'B' 'C']
```

```
Renamed DataFrame:
```

	ID	Group	Score_Value	Value_Rank
0	1	A	50	9.0
1	2	B	60	7.0
2	3	A	45	10.0
3	4	C	70	5.0
4	5	B	80	3.0
5	6	C	75	4.0
6	7	A	85	2.0
7	8	A	90	1.0
8	9	B	55	8.0
9	10	C	65	6.0

Week 6, 7, 8: Linear Regression and Residual Analysis

Develop a model on residual analysis of simple linear regression.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from scipy.stats import probplot

# Generate sample data for demonstration
np.random.seed(42)

X = np.random.rand(100, 1) * 10 # Predictor variable
y = 5 * X + np.random.normal(0, 2, (100, 1)) # Response variable with noise

# Split into training data
X = X.flatten() # Flatten for compatibility
y = y.flatten() # Flatten for compatibility

# 1. Develop a model on residual analysis of simple linear regression
model = LinearRegression()
model.fit(X.reshape(-1, 1), y)

# Predictions
y_pred = model.predict(X.reshape(-1, 1))
```

```
# Residuals
residuals = y - y_pred

print("predictions")
print(y_pred)
print("residuals")
print(residuals)

predictions
[18.81279122 47.09168167 36.35678697 29.81262224 8.08765795
8.08647412
 3.28096234 42.9425104 29.93318972 35.18271558 1.44048863
48.03380776
 41.28685468 10.85189217 9.35424417 9.43176884 15.36253959
26.18547598
 21.63024974 14.7238506 30.46020989 7.27661426 14.76878417
18.41139775
 22.81431399 38.96698461 10.2302721 25.66905178 29.50616804
2.70999964
 30.2487692 8.79959365 3.622954 47.00192485 47.82385046
40.10670081
 15.38077426 5.22398397 34.01265959 22.03307634 6.41988423
24.73369965
 2.11799653 45.06005007 13.13123118 32.94708652 15.7291133
25.95536667
 27.26298078 9.50293283 48.0178456 38.47406284 46.54122647
44.34872452
 29.77539445 45.67619788 4.77344445 10.04912009 2.64996817
16.3975525
 19.50664998 13.74812586 41.10500513 17.93980739 14.21858547
27.06596221
 7.34681719 39.80238396 4.08917172 48.86705076 38.33231584
10.18324809
 0.70122033 40.45340907 35.12307134 36.21019476 38.28449074
4.0643374
 18.02385285 6.11709831 42.79169991 31.02194707 16.67081745
3.54966497
 15.69334571 16.39033725 36.23959446 31.72180217 43.97499611
23.60671501
 6.29993215 35.43656996 37.76986755 27.97793174 38.26961118
24.66590414
 26.08615666 21.41409941 1.67777458 5.72555249]
residuals
[ 0.08830886 -0.15398105 0.42643167 -3.85483586 -0.7260697
0.42747704
 2.57900635 -0.67024355 -1.49442634 -0.78260077 1.41954032
1.11918707
 -0.72424304 0.79159823 -0.06884071 1.67574663 -1.55453363 -
0.60297869
 -0.81721511 -3.08942349 0.7246754 0.22018931 -0.15132483 -
```

```

0.56247985
-2.84155626 -0.54947719 -0.93201203 -1.5618844 -0.20801102
0.42062271
3.9008452 0.07576816 0.14472643 0.29346018 -3.37979123
0.26013885
-0.02962538 4.58610596 -0.18573019 0.57764303 -0.38739603 -
2.31221022
1.88707516 1.9098361 1.38983179 -1.63974721 2.66202912 -
2.75566774
1.24624737 4.1207012 -1.51968686 -0.85001714 0.63302334 -
0.61430831
-2.98172237 0.55463982 -2.47342678 0.69720789 -2.22745219
2.96883285
-1.63929208 -0.82479731 1.95890476 -2.56386969 0.28305989
2.68312745
-3.51557241 0.67673279 0.15812605 2.04094181 -2.1939788 -
2.88837724
0.61876866 0.9136317 0.72078155 0.93306006 -1.08102285
0.10240258
0.48557852 -1.75234817 4.0950204 1.09062511 -2.5085232
0.94135976
-2.09359296 1.44299806 2.55790562 -1.48529324 2.31239327
0.8295931
1.32390047 4.01925537 -0.22139135 -1.42154419 -1.50028104 -
1.6077449
-0.10371861 0.64525546 0.14656335 1.32338536]

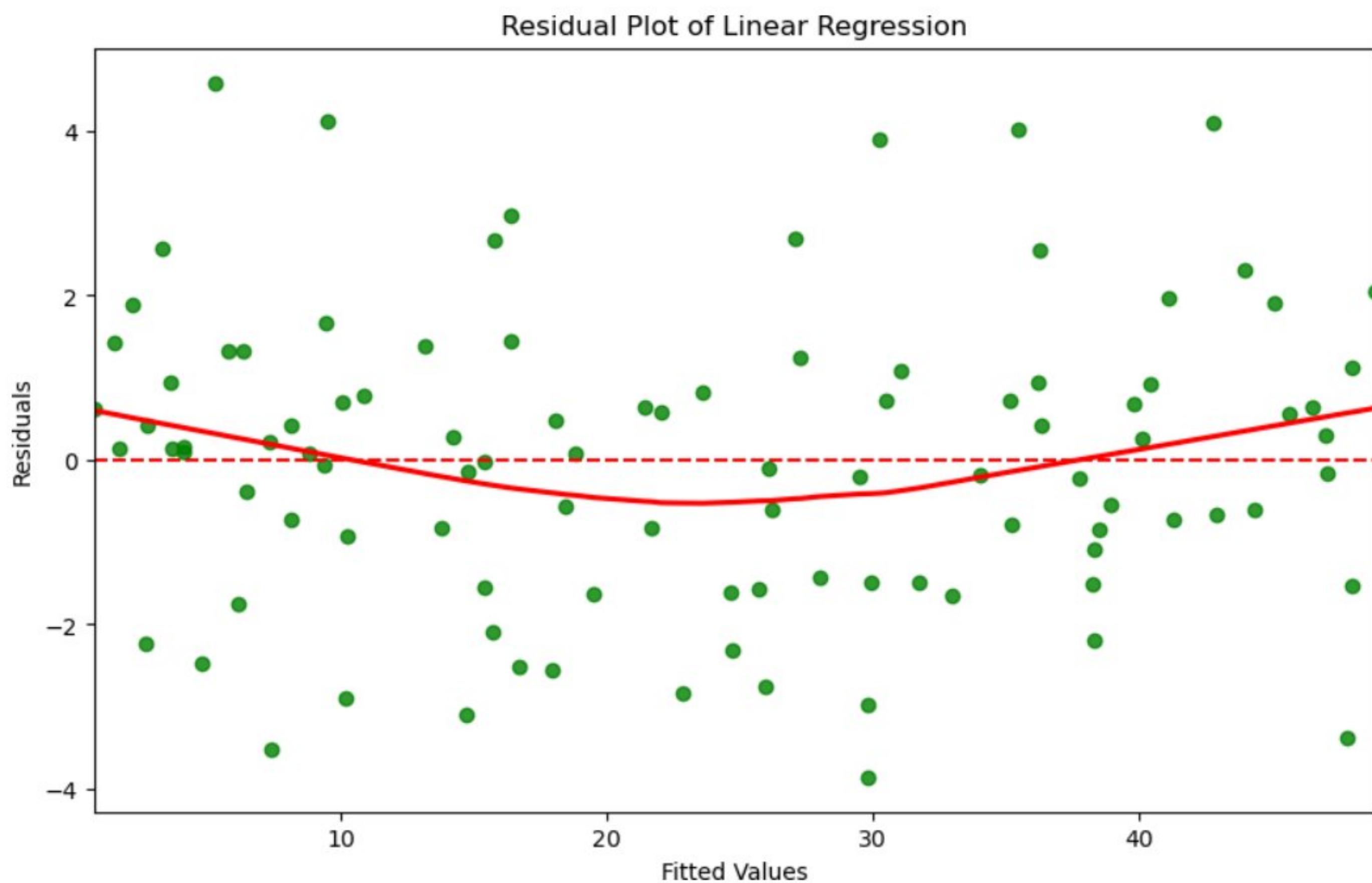
```

Residual plots of linear regression

```

plt.figure(figsize=(10, 6))
sns.residplot(x=y_pred, y=residuals, lowess=True, color='g',
line_kws={"color": "red"})
plt.axhline(y=0, color='r', linestyle='--')
plt.title("Residual Plot of Linear Regression")
plt.xlabel("Fitted Values")
plt.ylabel("Residuals")
plt.show()

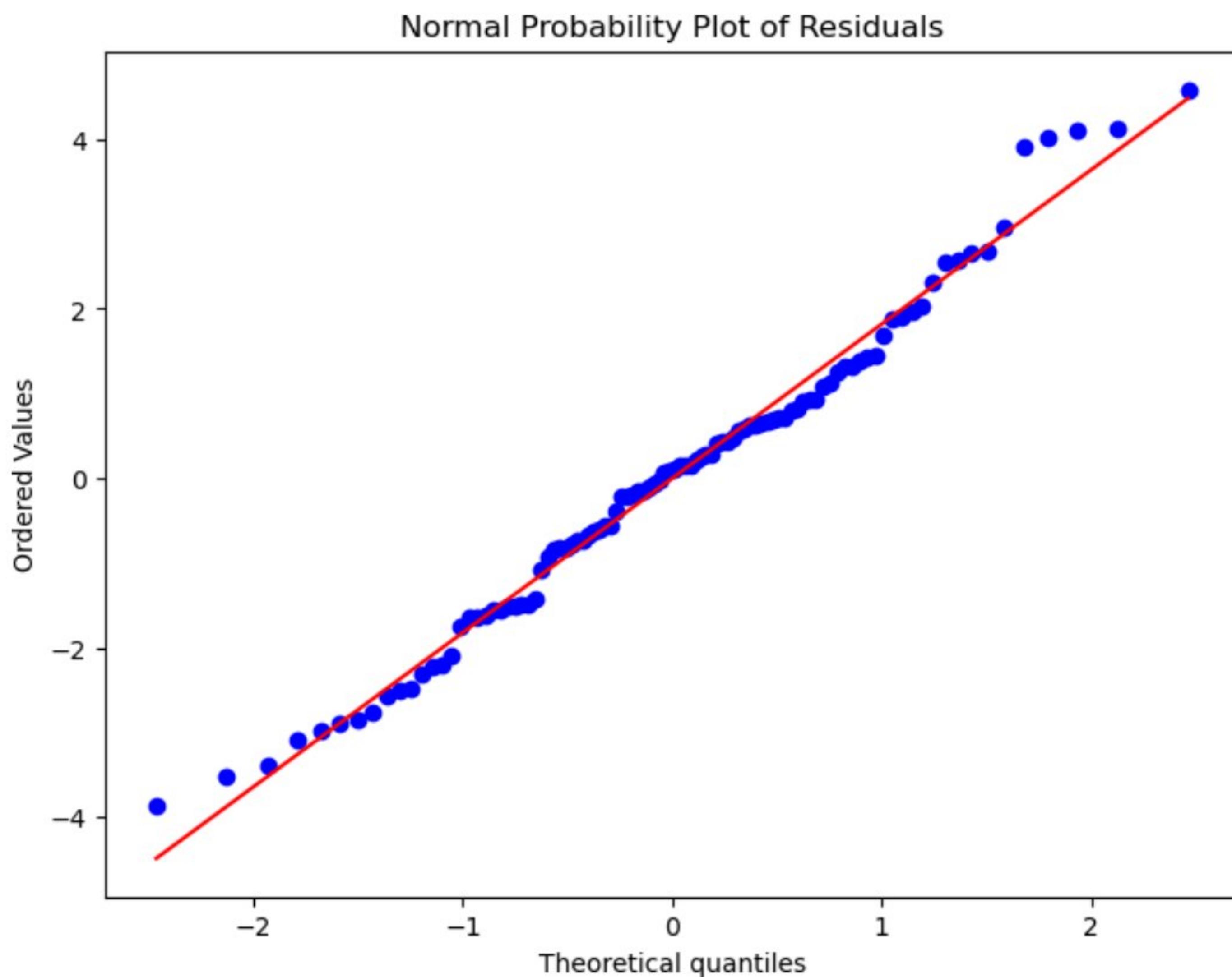
```



Normal probability plots

```
plt.figure(figsize=(8, 6))
probplot(residuals, dist="norm", plot=plt)

plt.title("Normal Probability Plot of Residuals")
plt.show()
```



Empirical model of linear regression analysis

```

# Print coefficients and intercept
print("Intercept:", model.intercept_)
print("Slope:", model.coef_[0])

# Empirical model equation
print(f"Empirical model: y = {model.intercept_:.2f} +
{model.coef_[0]:.2f} * X")

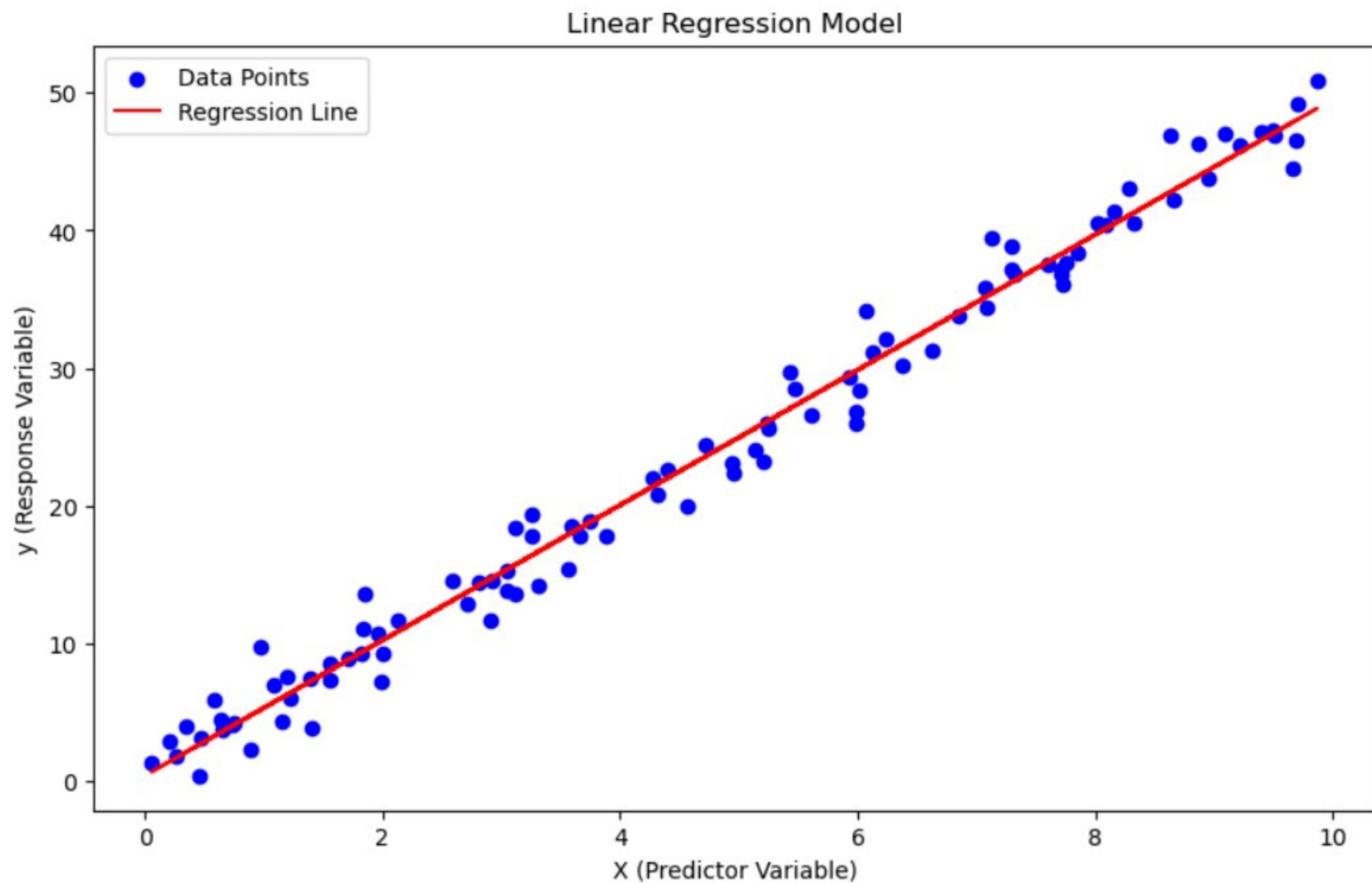
# Mean Squared Error
mse = mean_squared_error(y, y_pred)
print("Mean Squared Error:", mse)

# Visualize the regression line
plt.figure(figsize=(10, 6))
plt.scatter(X, y, color='blue', label='Data Points')
plt.plot(X, y_pred, color='red', label='Regression Line')
plt.title("Linear Regression Model")
plt.xlabel("X (Predictor Variable)")

```

```
plt.ylabel("y (Response Variable)")  
plt.legend()  
plt.show()
```

```
Intercept: 0.43019231509349254  
Slope: 4.9080453545753935  
Empirical model: y = 0.43 + 4.91 * x  
Mean Squared Error: 3.226338255868214
```



Week 9, 10, 11 Tasks

Task 1: Pandas DataFrame Operations

```
#Step 1: Importing a CSV File
import pandas as pd
# Import CSV file
df = pd.read_csv('iris.csv')
# Display the first few rows of the DataFrame
print(df.head())

  sepal_length  sepal_width  petal_length  petal_width species
0          5.1         3.5          1.4         0.2   setosa
1          4.9         3.0          1.4         0.2   setosa
```

```

2      4.7      3.2      1.3      0.2  setosa
3      4.6      3.1      1.5      0.2  setosa
4      5.0      3.6      1.4      0.2  setosa

#Step 2: Handling Missing Data
#(a) Detecting and Dropping/ Filling Missing Values
# Detect missing values
print(df.isnull().sum())

# Drop rows with missing values
df_dropped = df.dropna()
print(df_dropped)

# Fill missing values with a specific value or method
df_filled = df.fillna(method='ffill') # Forward fill
print(df_filled)

sepal_length    0
sepal_width     0
petal_length    0
petal_width     0
species         0
dtype: int64
   sepal_length  sepal_width  petal_length  petal_width  species
0            5.1        3.5          1.4        0.2  setosa
1            4.9        3.0          1.4        0.2  setosa
2            4.7        3.2          1.3        0.2  setosa
3            4.6        3.1          1.5        0.2  setosa
4            5.0        3.6          1.4        0.2  setosa
..
145           6.7        3.0          5.2        2.3  virginica
146           6.3        2.5          5.0        1.9  virginica
147           6.5        3.0          5.2        2.0  virginica
148           6.2        3.4          5.4        2.3  virginica
149           5.9        3.0          5.1        1.8  virginica

[150 rows x 5 columns]
   sepal_length  sepal_width  petal_length  petal_width  species
0            5.1        3.5          1.4        0.2  setosa
1            4.9        3.0          1.4        0.2  setosa
2            4.7        3.2          1.3        0.2  setosa
3            4.6        3.1          1.5        0.2  setosa
4            5.0        3.6          1.4        0.2  setosa
..
145           6.7        3.0          5.2        2.3  virginica
146           6.3        2.5          5.0        1.9  virginica
147           6.5        3.0          5.2        2.0  virginica
148           6.2        3.4          5.4        2.3  virginica
149           5.9        3.0          5.1        1.8  virginica

```

```
[150 rows x 5 columns]

C:\Users\B11202016\AppData\Local\Temp\ipykernel_5692\3643678138.py:11:
FutureWarning: DataFrame.fillna with 'method' is deprecated and will
raise in a future version. Use obj.ffill() or obj.bfill() instead.
    df_filled = df.fillna(method='ffill') # Forward fill

#Step 3: Data Transformation Using apply() and map()
# Using apply() for row-wise or column-wise operations
# Transform data using apply () and map() method.

import pandas as pd
# Sample DataFrame
data = {
    'A': [1, 2, 3],
    'B': [10, 20, 30],
    'C': [100, 200, 300]
}
df = pd.DataFrame(data)

# Define a function
def multiply_by_2(x):
    return x * 2

# Apply the function to the DataFrame
df_applied = df.apply(multiply_by_2)
print("DataFrame after apply():")
print(df_applied)

DataFrame after apply():
   A   B   C
0  2  20  200
1  4  40  400
2  6  60  600

data = {
    'A': [1, 2, 3],
    'B': [10, 20, 30],
    'C': [100, 200, 300]
}
df = pd.DataFrame(data)

df['A'] = df['A'].map(multiply_by_2)
print(df)

   A   B   C
0  2  10  100
1  4  20  200
2  6  30  300
```

(c) Detect and filter outliers.

```
import pandas as pd

# Sample DataFrame
data = {
    'A': [1, 2, 3, 4, 5, 100, 6, 7, 8, 9]
}
df = pd.DataFrame(data)

# Function to detect and filter out outliers using IQR
def remove_outliers(df, column_name):
    Q1 = df[column_name].quantile(0.25)
    Q3 = df[column_name].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

    # Filter out outliers
    filtered_df = df[(df[column_name] >= lower_bound) &
                      (df[column_name] <= upper_bound)]

    return filtered_df

# Detect and filter out outliers from column 'A'
filtered_df = remove_outliers(df, 'A')
print("Filtered DataFrame (outliers removed):")
print(filtered_df)

Filtered DataFrame (outliers removed):
   A
0  1
1  2
2  3
3  4
4  5
6  6
7  7
8  8
9  9

import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Generate a simple linear dataset
np.random.seed(42)
X = np.random.rand(50, 1) * 10  # Random values between 0 and 10
```

```

y = 3 * X + np.random.randn(50, 1) * 3 # y = 3x + noise

# Add outliers
X_outliers = np.append(X, [[20], [22]], axis=0) # Adding extreme X values
y_outliers = np.append(y, [[50], [60]], axis=0) # Adding extreme Y values

# Create Linear Regression models
model_no_outliers = LinearRegression()
model_with_outliers = LinearRegression()

# Fit models
model_no_outliers.fit(X, y)
model_with_outliers.fit(X_outliers, y_outliers)

# Predictions
y_pred_no_outliers = model_no_outliers.predict(X)
y_pred_with_outliers = model_with_outliers.predict(X_outliers)

# Mean Squared Error (MSE)
mse_no_outliers = mean_squared_error(y, y_pred_no_outliers)
mse_with_outliers = mean_squared_error(y_outliers,
y_pred_with_outliers)

# Plotting
plt.figure(figsize=(14, 6))

# Without Outliers
plt.subplot(1, 2, 1)
plt.scatter(X, y, color='blue', label='Data Points')
plt.plot(X, y_pred_no_outliers, color='red', label='Regression Line')
plt.title("Without Outliers")
plt.xlabel("X")
plt.ylabel("y")
plt.legend()
plt.grid(True)

# With Outliers
plt.subplot(1, 2, 2)
plt.scatter(X_outliers, y_outliers, color='blue', label='Data Points')
plt.plot(X_outliers, y_pred_with_outliers, color='red',
label='Regression Line')
plt.title("With Outliers")
plt.xlabel("X")
plt.ylabel("y")
plt.legend()
plt.grid(True)

plt.tight_layout()

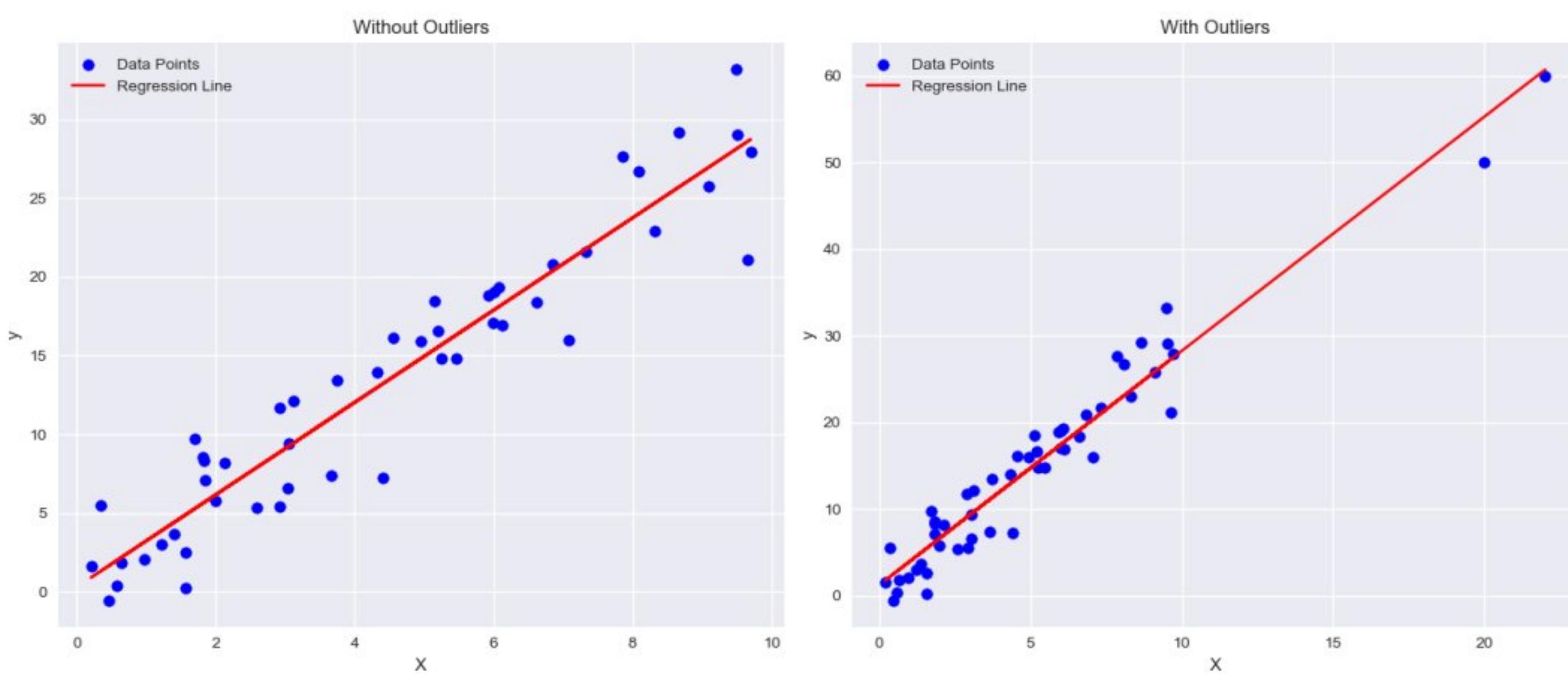
```

```

plt.show()

# Print Mean Squared Errors
print(f"Mean Squared Error (Without Outliers): {mse_no_outliers:.2f}")
print(f"Mean Squared Error (With Outliers): {mse_with_outliers:.2f}")

```



```

Mean Squared Error (Without Outliers): 7.41
Mean Squared Error (With Outliers): 8.09

```

(d) Perform Vectorized String operations on Pandas Series. 2. Implement regularized Linear regression

```

import pandas as pd

# Sample data with strings
data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eva']}
df = pd.DataFrame(data)

# Convert strings to uppercase
df['Name_upper'] = df['Name'].str.upper()

# Convert strings to lowercase
df['Name_lower'] = df['Name'].str.lower()

# Replace a substring (e.g., replace 'a' with 'X')
df['Name_replaced'] = df['Name'].str.replace('a', 'X')

# Check if the string contains the letter 'e'
df['Has_e'] = df['Name'].str.contains('e')

# Extract the first 3 characters of each name
df['Name_first3'] = df['Name'].str[:3]

```

```

# Count the occurrences of the letter 'a' in each name
df['Count_a'] = df['Name'].str.count('a')

print("DataFrame after Vectorized String Operations:")
print(df)

DataFrame after Vectorized String Operations:
   Name  Name_upper  Name_lower  Name_replaced  Has_e  Name_first3
Count_a
0    Alice      ALICE      alice        Alice  True    Ali
0
1     Bob       BOB       bob        Bob  False   Bob
0
2  Charlie    CHARLIE    charlie    ChXrlie  True   Cha
1
3   David     DAVID     david     DXvid  False   Dav
1
4    Eva       EVA       eva      EvX  False   Eva
1

```

e) Develop a model on logistic regression on any data set for prediction

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.datasets import load_iris

# Load the Iris dataset from sklearn
iris = load_iris()
X = iris.data # Features (sepal_length, sepal_width, petal_length,
petal_width)
y = iris.target # Target labels (species)

# Convert to DataFrame for easier handling
df = pd.DataFrame(X, columns=iris.feature_names)
df['species'] = iris.target

# Display the first few rows
print("First 5 rows of the dataset:")
print(df.head())

# Split the data into training and test sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

# Standardize the features (important for logistic regression)
scaler = StandardScaler()

```

```

X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Initialize the Logistic Regression model
logreg = LogisticRegression(multi_class='ovr', max_iter=200)

# Train the model
logreg.fit(X_train_scaled, y_train)

# Make predictions on the test set
y_pred = logreg.predict(X_test_scaled)

# Model Evaluation
print("\nClassification Report:")
print(classification_report(y_test, y_pred))

print("\nConfusion Matrix:")
print(confusion_matrix(y_test, y_pred))

# Predict on new data (example)
new_data = np.array([[5.4, 3.9, 1.7, 0.4], # Example flower 1
                     [6.7, 3.1, 4.7, 1.5]]) # Example flower 2

new_data_scaled = scaler.transform(new_data)
predictions = logreg.predict(new_data_scaled)

# Print the species predictions (0 -> Setosa, 1 -> Versicolor, 2 ->
# Virginica)
print("\nPredictions for new data:")
print(predictions)

First 5 rows of the dataset:
   sepal length (cm)  sepal width (cm)  petal length (cm)  petal width
   (cm) \
0           5.1              3.5            1.4
0.2
1           4.9              3.0            1.4
0.2
2           4.7              3.2            1.3
0.2
3           4.6              3.1            1.5
0.2
4           5.0              3.6            1.4
0.2

   species
0       0
1       0
2       0
3       0

```

```
4      0
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	10
1	1.00	0.89	0.94	9
2	0.92	1.00	0.96	11
accuracy			0.97	30
macro avg	0.97	0.96	0.97	30
weighted avg	0.97	0.97	0.97	30

Confusion Matrix:

```
[ 0  0 11]]
```

Predictions for new data:

```
[0 1]
```



1. Visualize Data Using Line Plots, Bar Plots, Histograms, Density Plots, and Scatter Plots

```
In [ ]: import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
tips = sns.load_dataset('tips')

# Line Plot: Total bill vs Tip
plt.figure(figsize=(12, 8))
sns.lineplot(data=tips, x='total_bill', y='tip', hue='sex')
plt.title('Line Plot: Tip vs Total Bill by Sex')
plt.show()

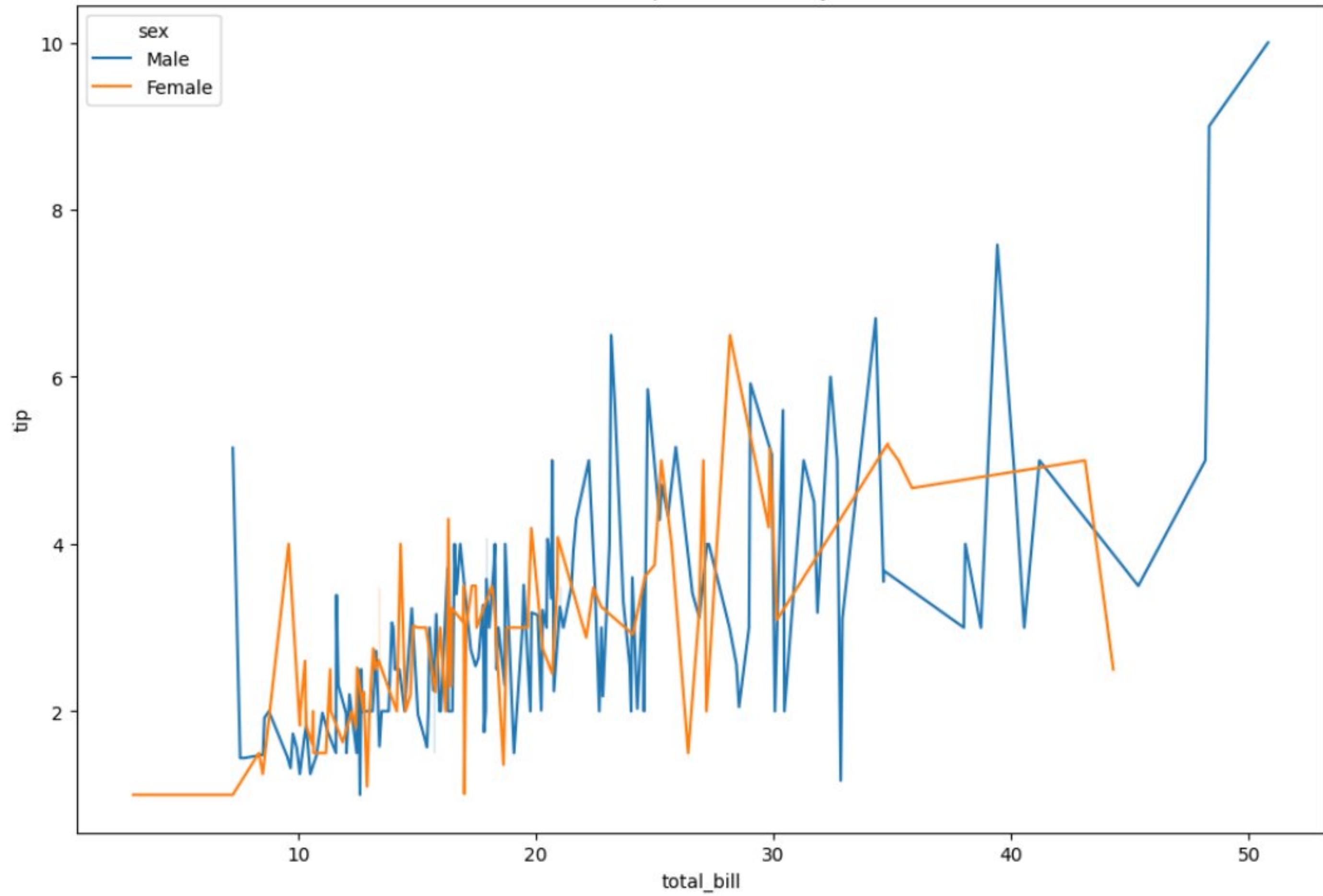
# Bar Plot: Average tip by day
plt.figure(figsize=(12, 8))
sns.barplot(data=tips, x='day', y='tip', hue='time')
plt.title('Bar Plot: Average Tip by Day and Time')
plt.show()

# Histogram: Distribution of total bill
plt.figure(figsize=(12, 8))
sns.histplot(data=tips, x='total_bill', bins=20, kde=True)
plt.title('Histogram: Distribution of Total Bill')
plt.show()

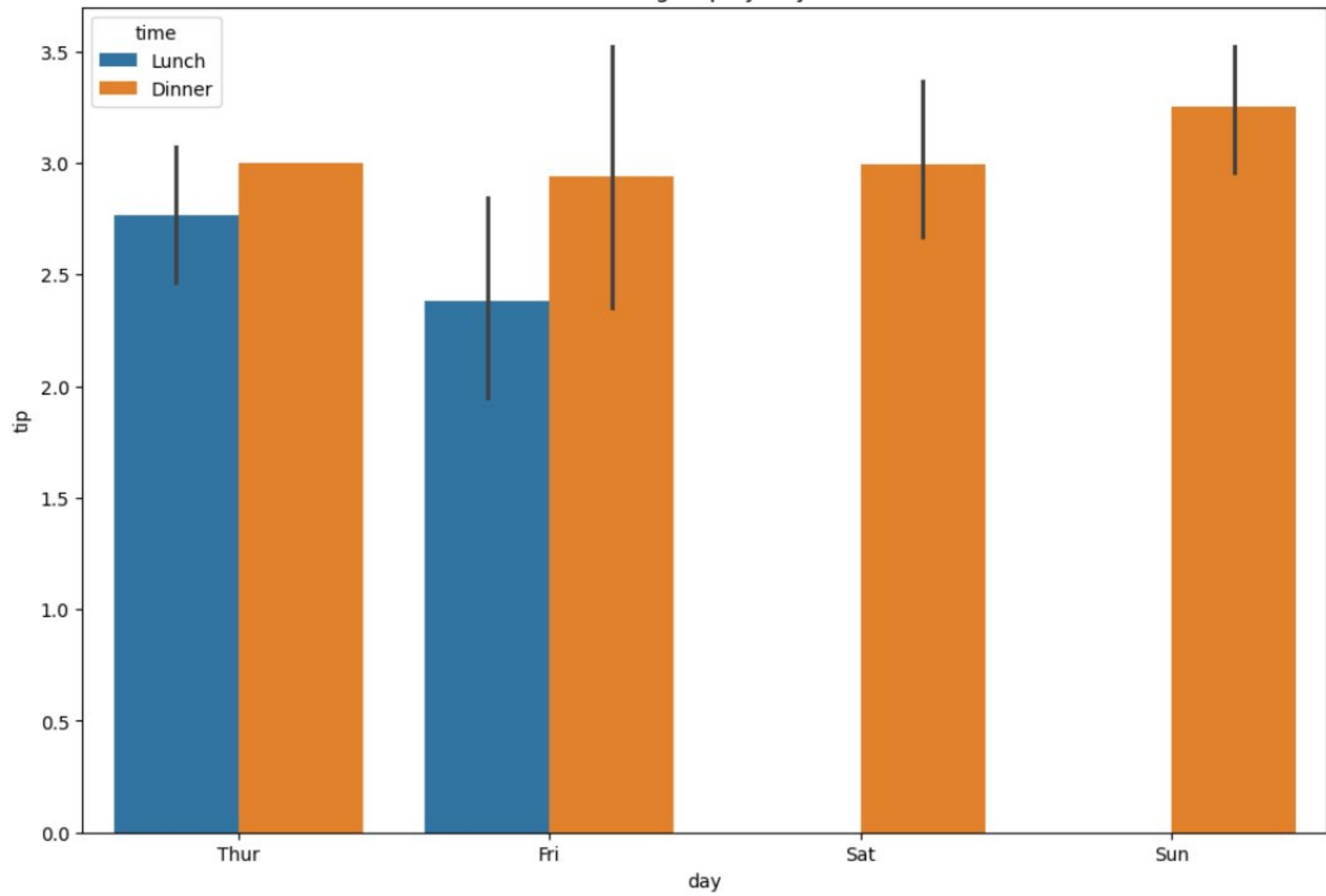
# Density Plot: Density of tip amounts
plt.figure(figsize=(12, 8))
sns.kdeplot(data=tips, x='tip', hue='smoker', fill=True)
plt.title('Density Plot: Tip Amounts by Smoker Status')
plt.show()

# Scatter Plot: Total bill vs Tip
plt.figure(figsize=(12, 8))
sns.scatterplot(data=tips, x='total_bill', y='tip', hue='size', size='size')
plt.title('Scatter Plot: Total Bill vs Tip (Size by Party Size)')
plt.show()
```

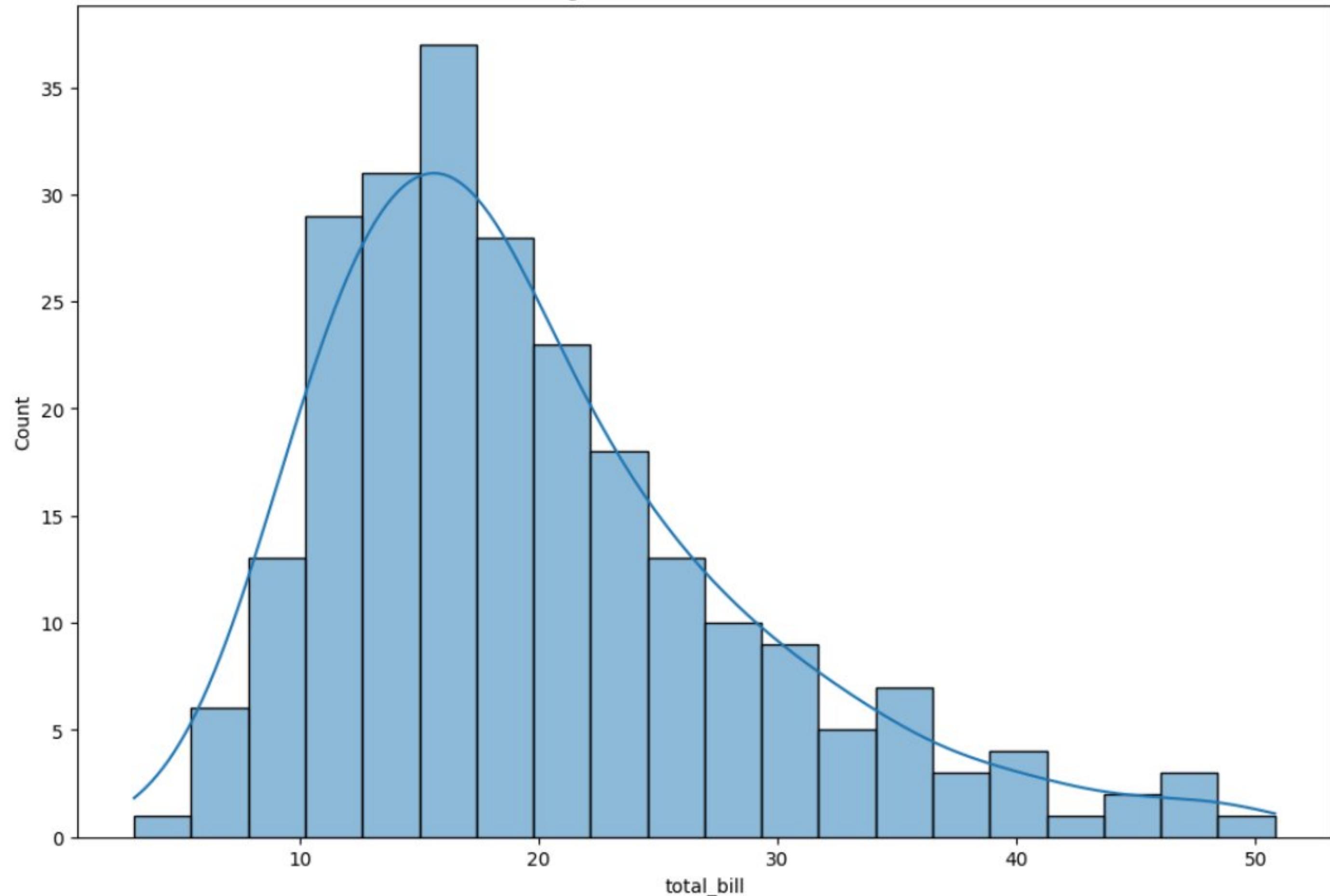
Line Plot: Tip vs Total Bill by Sex



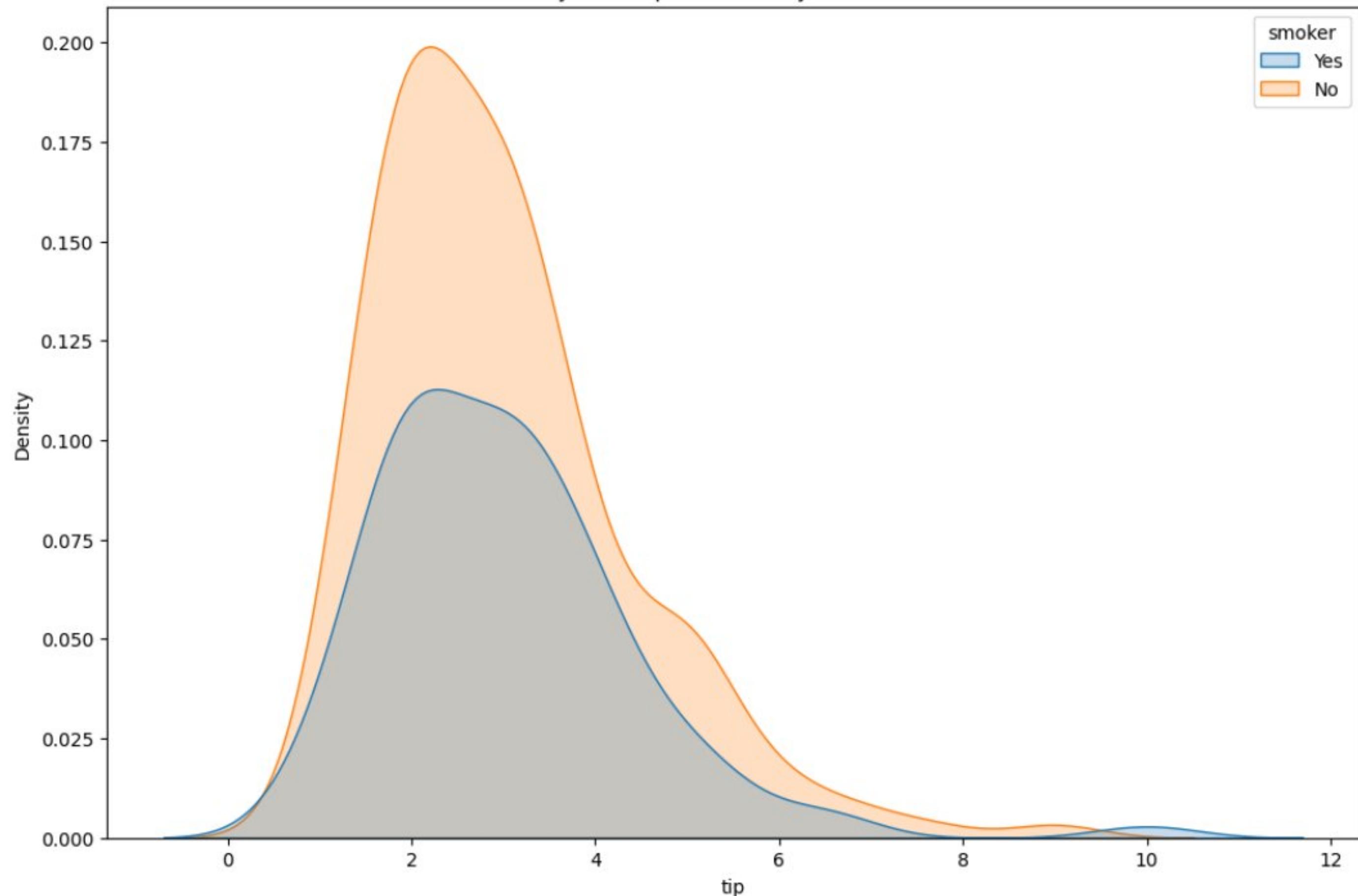
Bar Plot: Average Tip by Day and Time

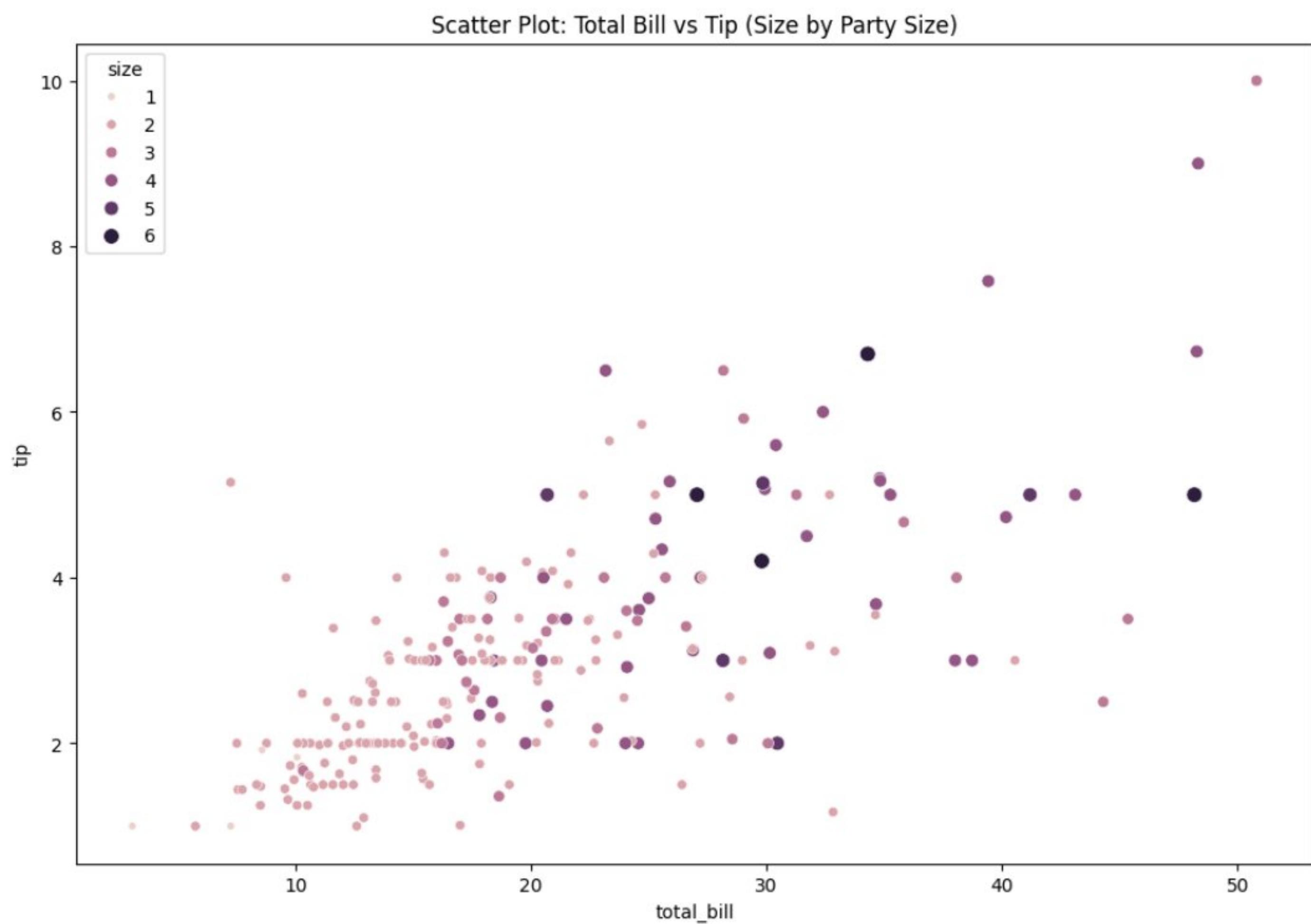


Histogram: Distribution of Total Bill



Density Plot: Tip Amounts by Smoker Status





2. Download the House Pricing Dataset from Kaggle and Map Values to 23 Aesthetics

```
In [ ]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset (assume downloaded as 'train.csv' in current directory)
df = pd.read_csv('train.csv')

# Select 23 key features for aesthetics mapping (mix of categorical/numerical)
aesthetic_features = [
    'MSZoning', 'LotShape', 'LandContour', 'Utilities', 'LotConfig', 'LandSlope',
    'Neighborhood', 'Condition1', 'Condition2', 'BldgType', 'HouseStyle',
    'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd', 'RoofStyle',
    'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType', 'ExterQual',
    'Foundation', 'BsmtQual'
]

# For demonstration, we'll cycle through the first 10 for hue/style/size (to a
# Full 23 can be extended by rotating in a loop if needed
subset_aesthetics = aesthetic_features[:10]

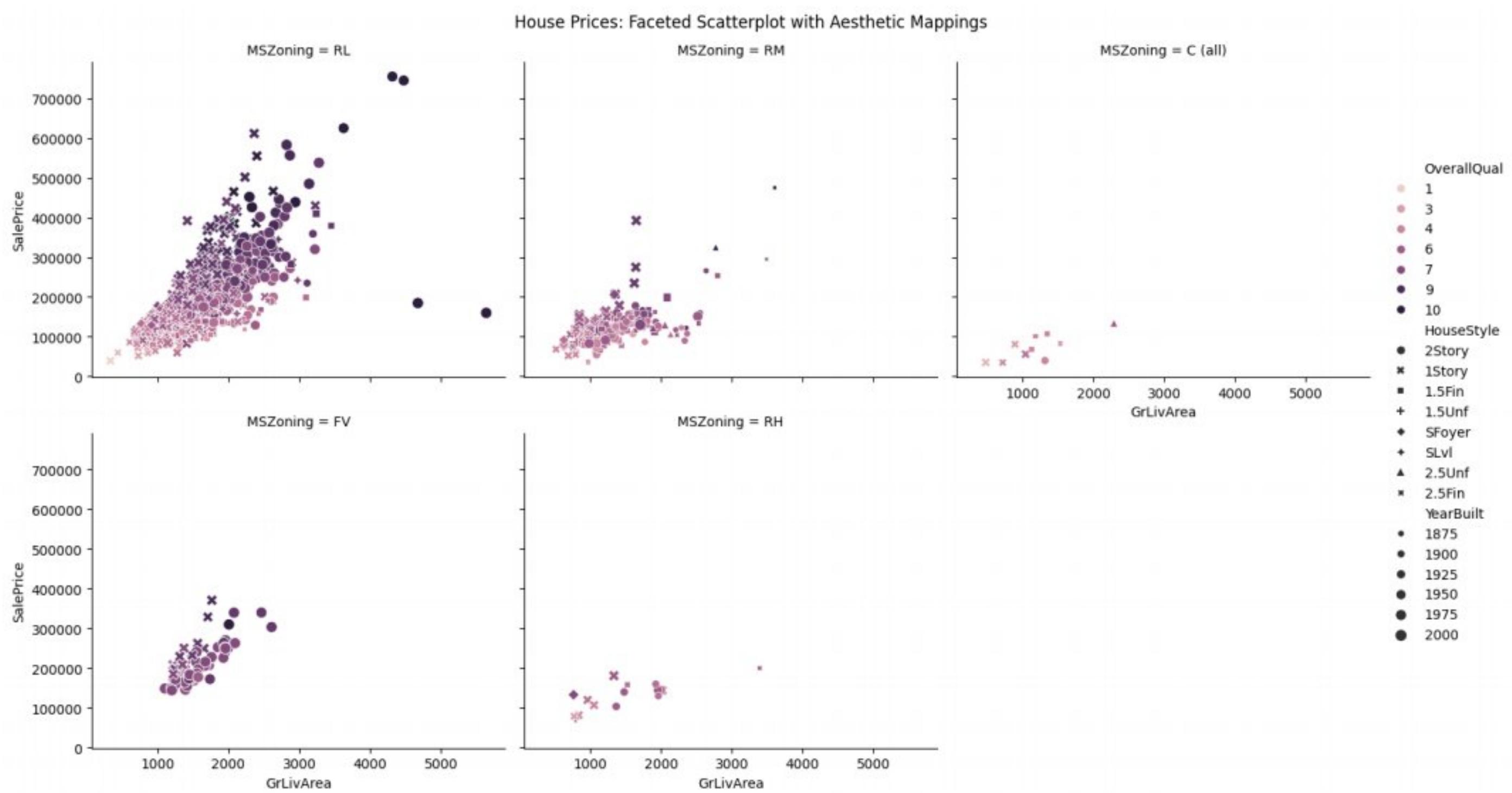
# Create a faceted scatterplot mapping aesthetics
g = sns.relplot(
```

```

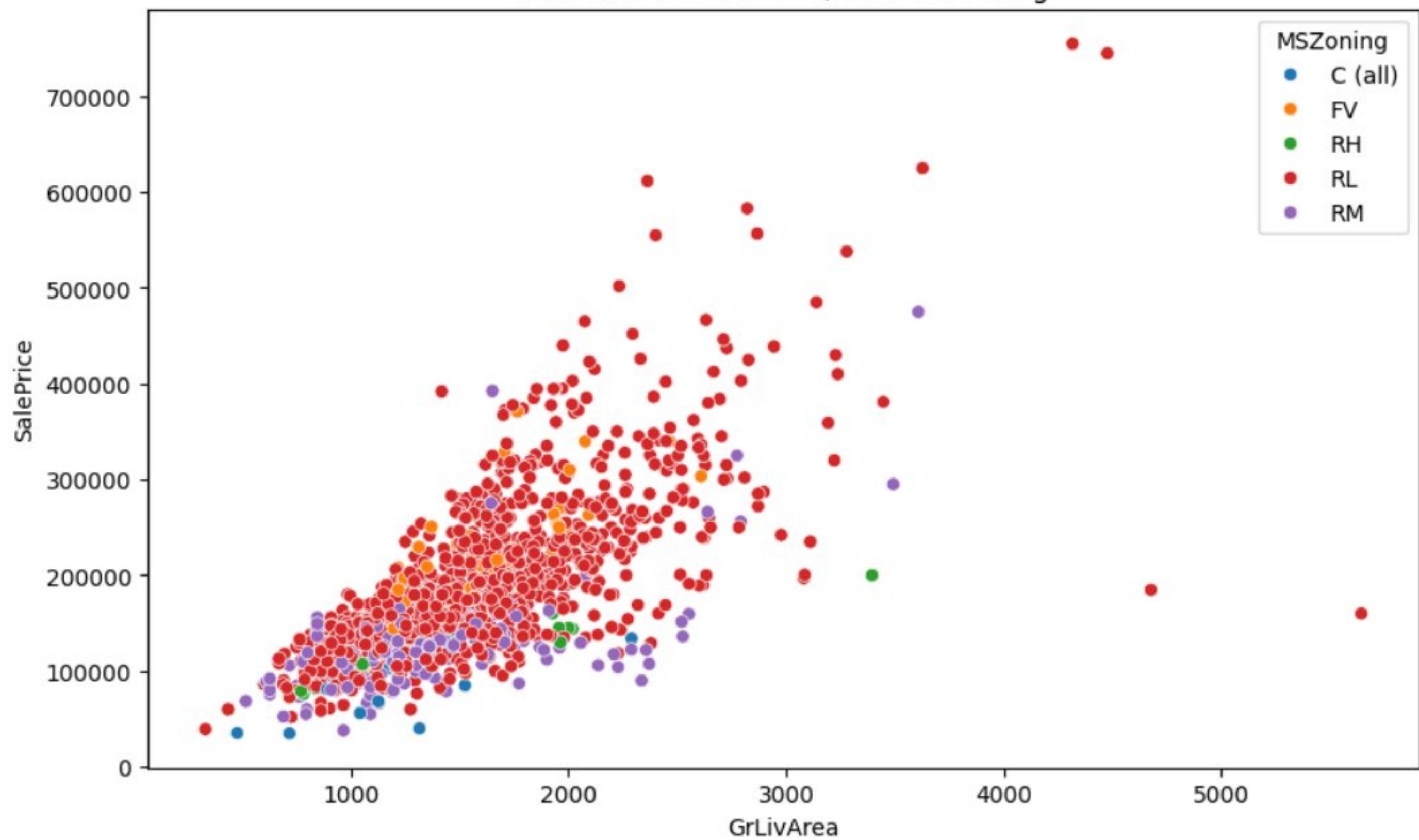
    data=df, x='GrLivArea', y='SalePrice',
    col='MSZoning', col_wrap=3, # Facet by one feature
    hue='OverallQual', # Color by quality
    size='YearBuilt', # Size by year (normalized)
    style='HouseStyle', # Style by house style
    height=4, aspect=1.2
)
plt.suptitle('House Prices: Faceted Scatterplot with Aesthetic Mappings', y=1.
plt.show()

# To map all 23: Loop and print individual plots (example for first 5)
for i, feat in enumerate(aesthetic_features[:5]):
    plt.figure(figsize=(10, 6))
    sns.scatterplot(data=df, x='GrLivArea', y='SalePrice', hue=df[feat].astype
    plt.title(f'SalePrice vs GrLivArea, Hue: {feat}')
    plt.show()

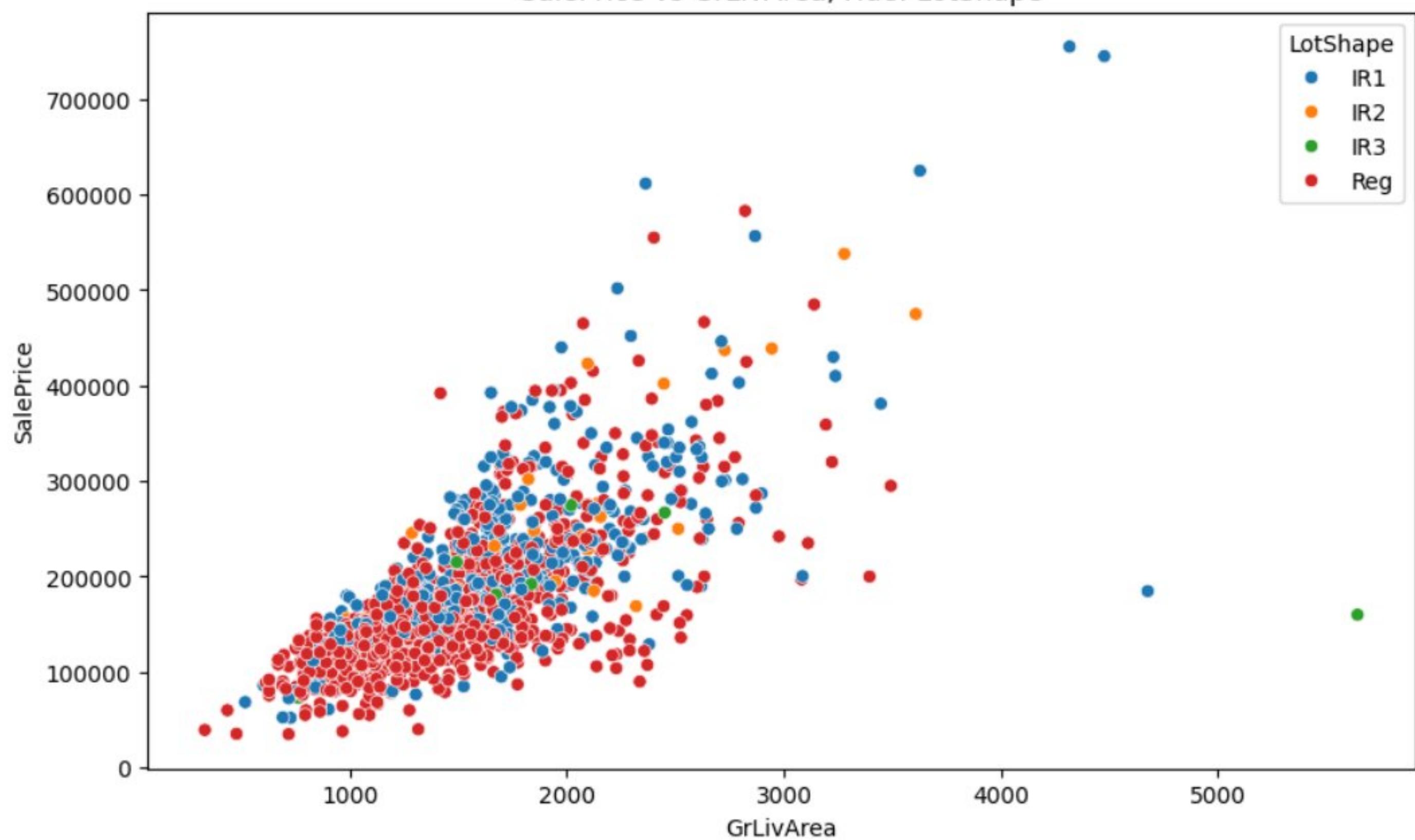
```



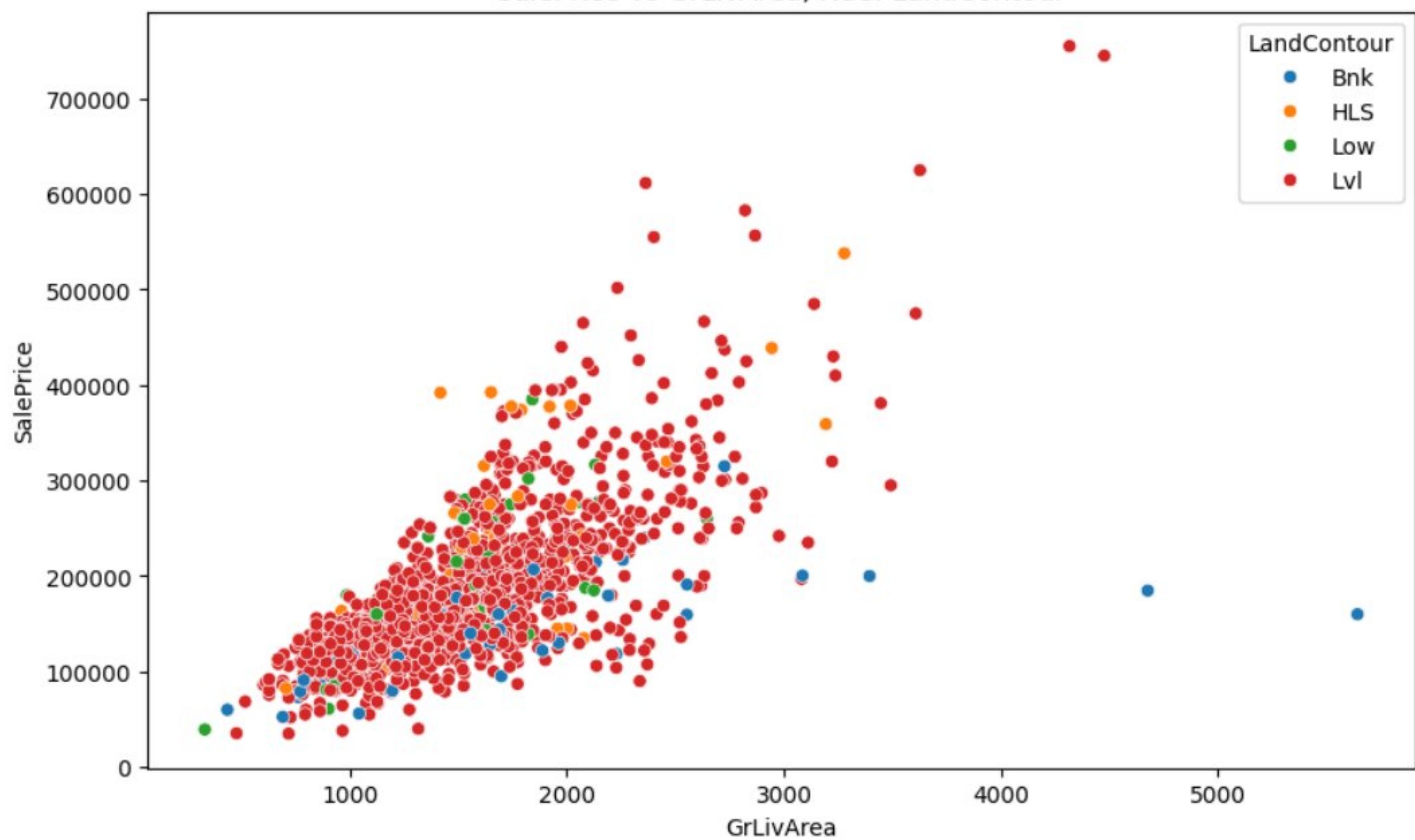
SalePrice vs GrLivArea, Hue: MSZoning



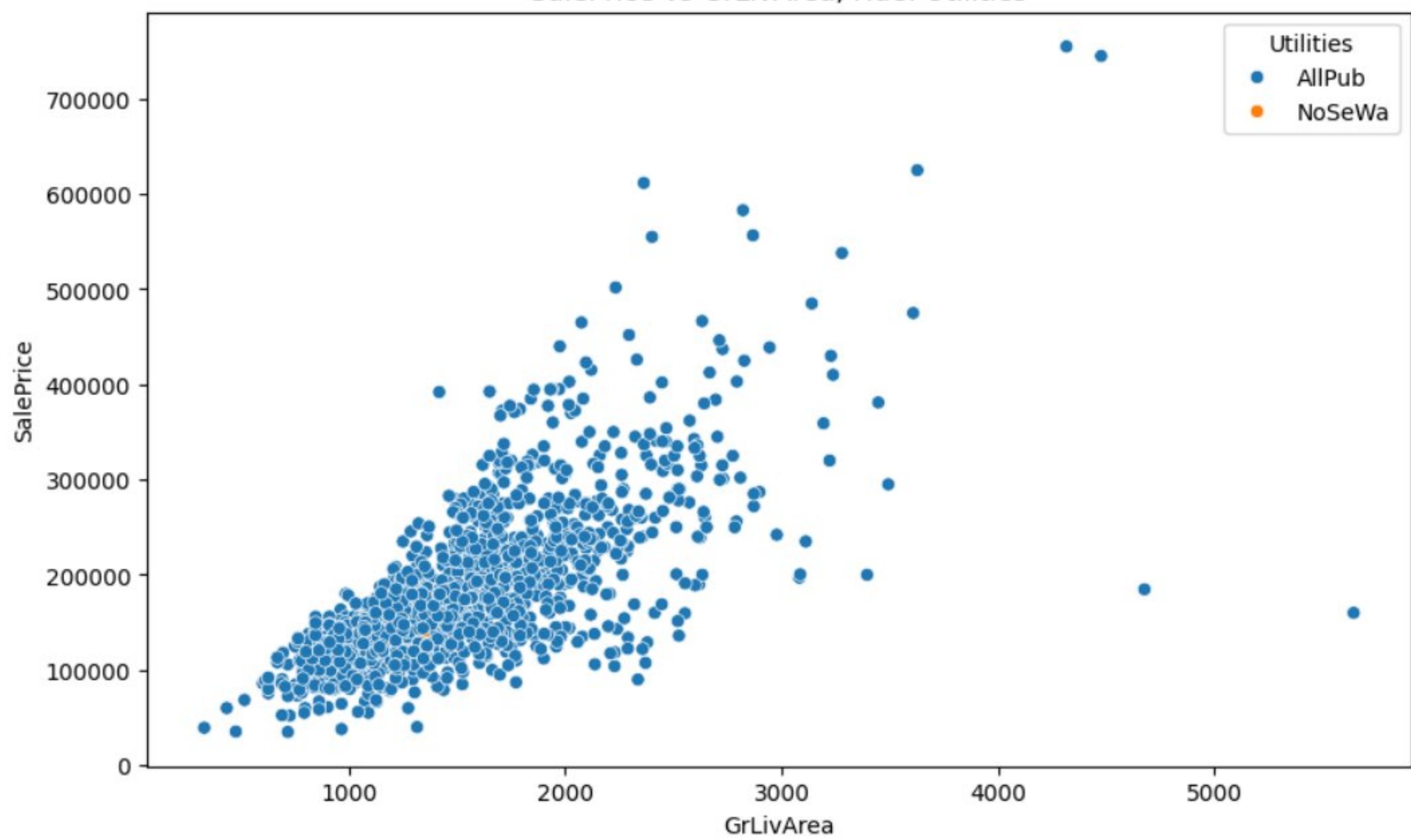
SalePrice vs GrLivArea, Hue: LotShape

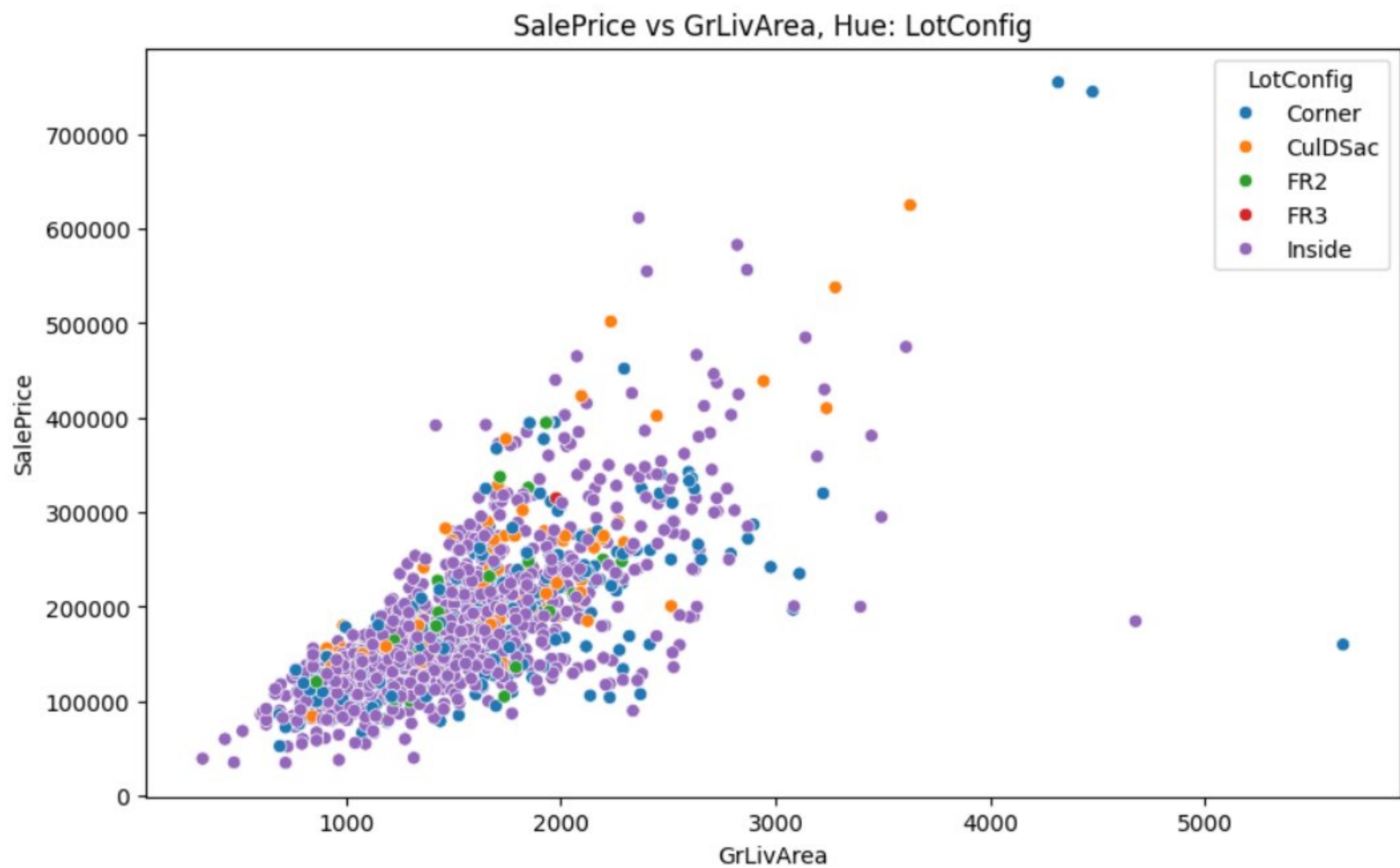


SalePrice vs GrLivArea, Hue: LandContour



SalePrice vs GrLivArea, Hue: Utilities





4. Create Different Bar Plots for Variables in Any Dataset

```
In [ ]: import seaborn as sns
import matplotlib.pyplot as plt

# Load dataset
flights = sns.load_dataset('flights')

# Horizontal Bar Plot: Passengers by Month (averaged over years)
monthly_avg = flights.groupby('month').mean().reset_index()
plt.figure(figsize=(10, 6))
sns.barplot(data=monthly_avg, y='month', x='passengers', orient='h')
plt.title('Horizontal Bar: Average Passengers by Month')
plt.show()

# Grouped Bar Plot: Passengers by Year and Month
plt.figure(figsize=(12, 6))
sns.barplot(data=flights, x='year', y='passengers', hue='month')
plt.title('Grouped Bar: Passengers by Year and Month')
plt.show()

# Stacked Bar Plot (using pandas plot for stacking)
yearly_monthly = flights.set_index(['year', 'month'])['passengers'].unstack()
plt.figure(figsize=(12, 6))
yearly_monthly.plot(kind='bar', stacked=True, figsize=(12, 6))
plt.title('Stacked Bar: Passengers by Year (Stacked by Month)')
plt.xlabel('Year')
plt.ylabel('Passengers')
```

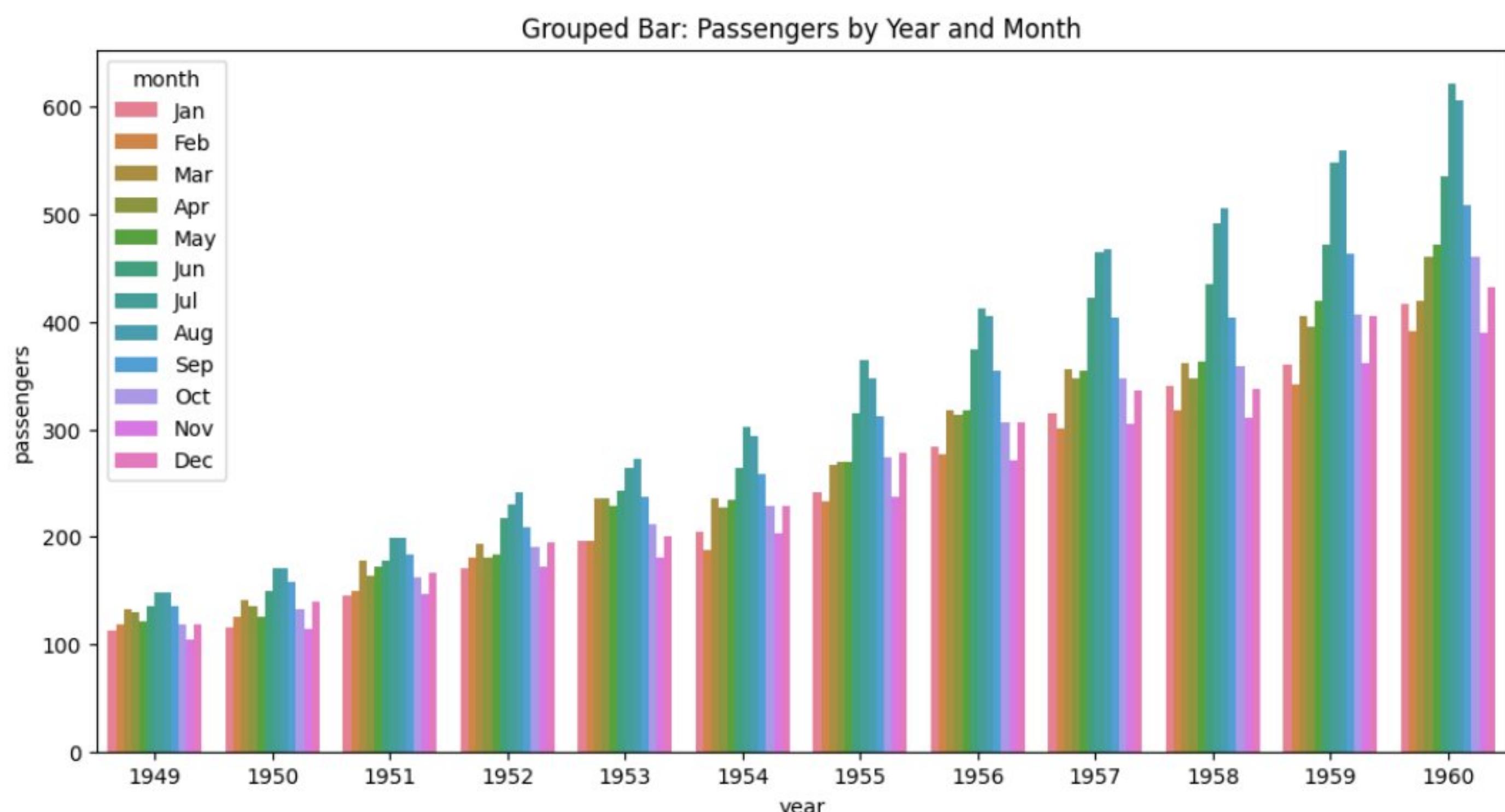
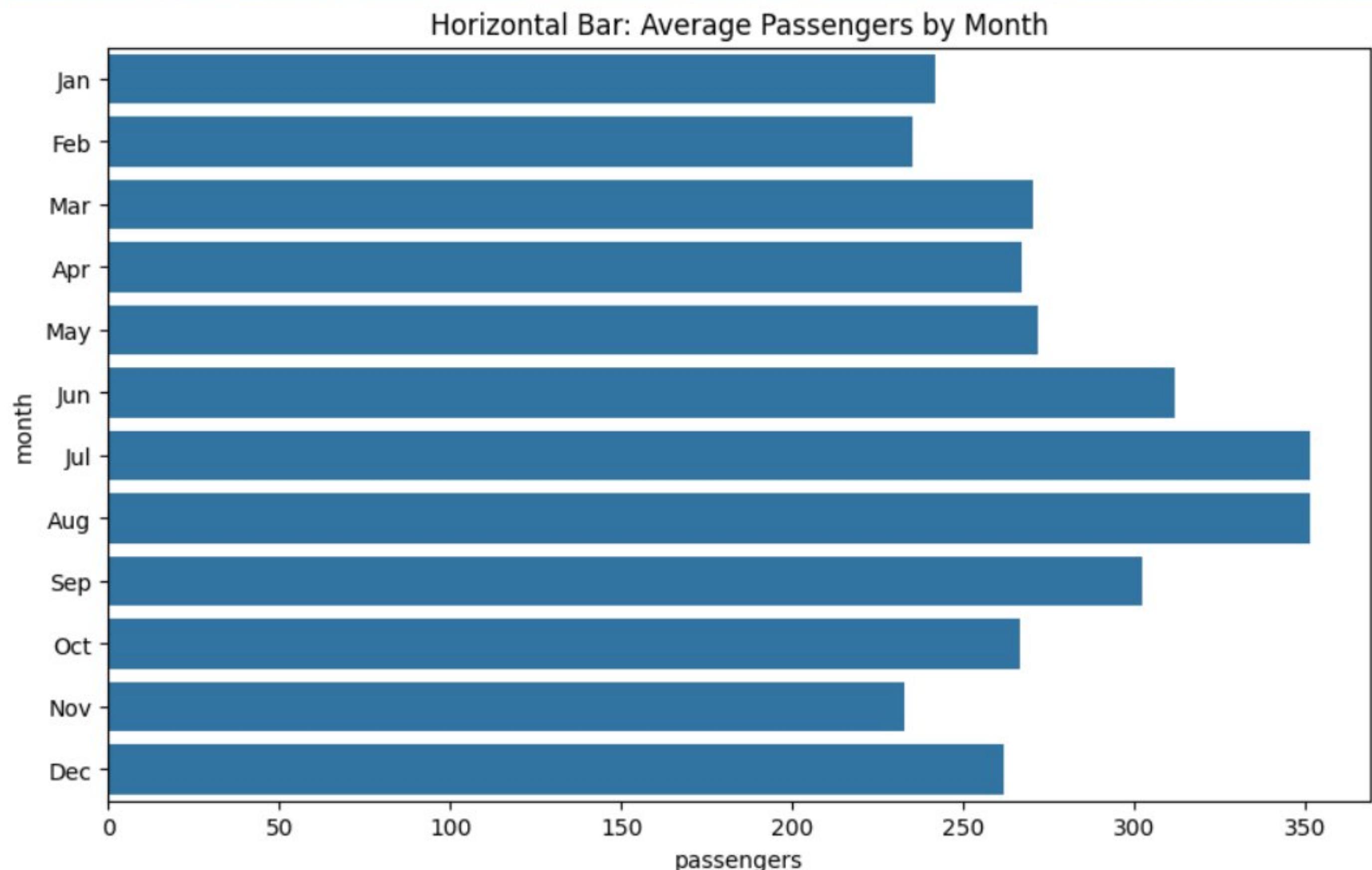
```

plt.legend(title='Month', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()

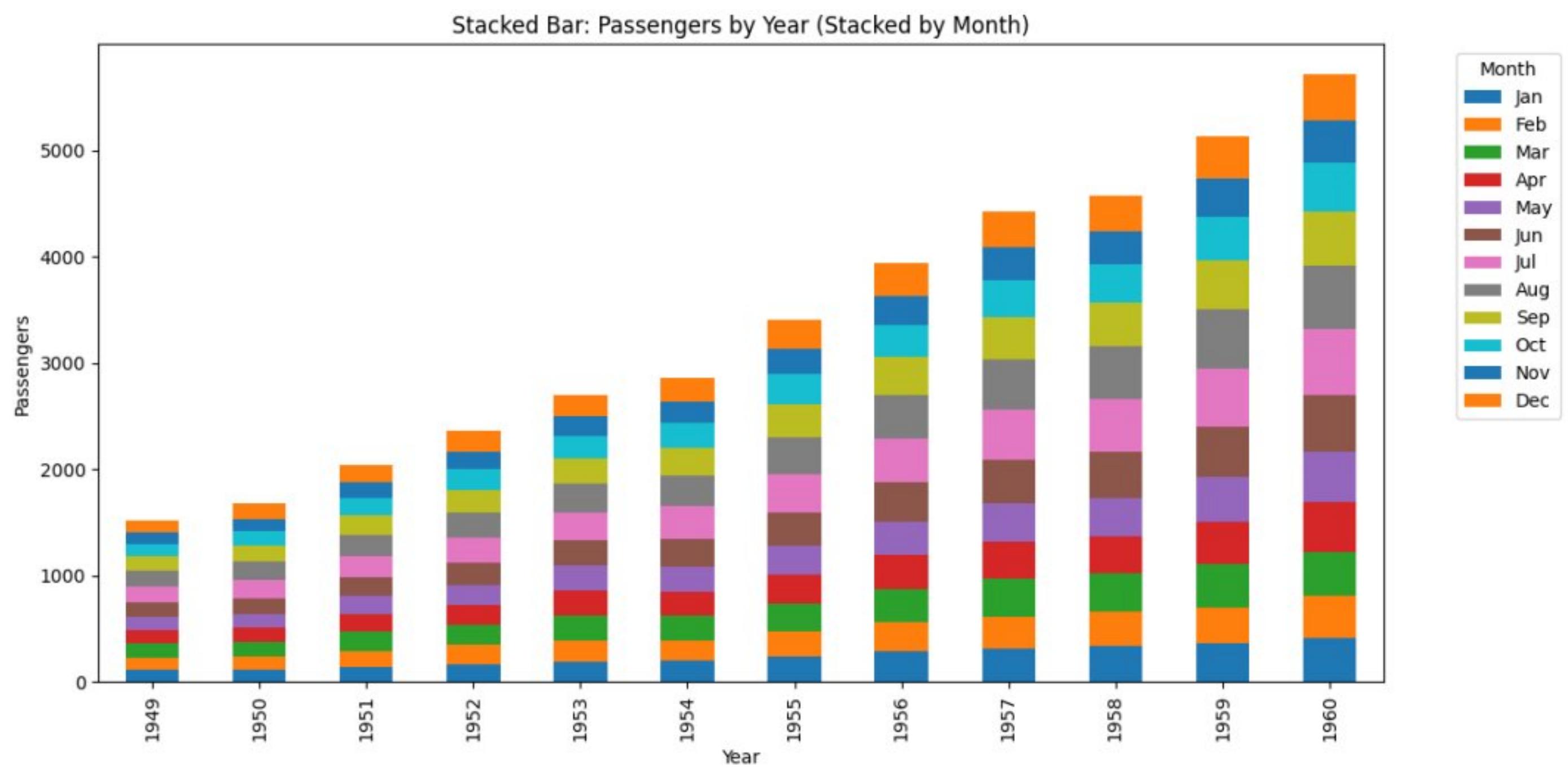
```

/tmp/ipython-input-805159941.py:8: 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.

```
monthly_avg = flights.groupby('month').mean().reset_index()
```



<Figure size 1200x600 with 0 Axes>



5. Show an Example of Skewed Data and Removal of Skewness

```
In [ ]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

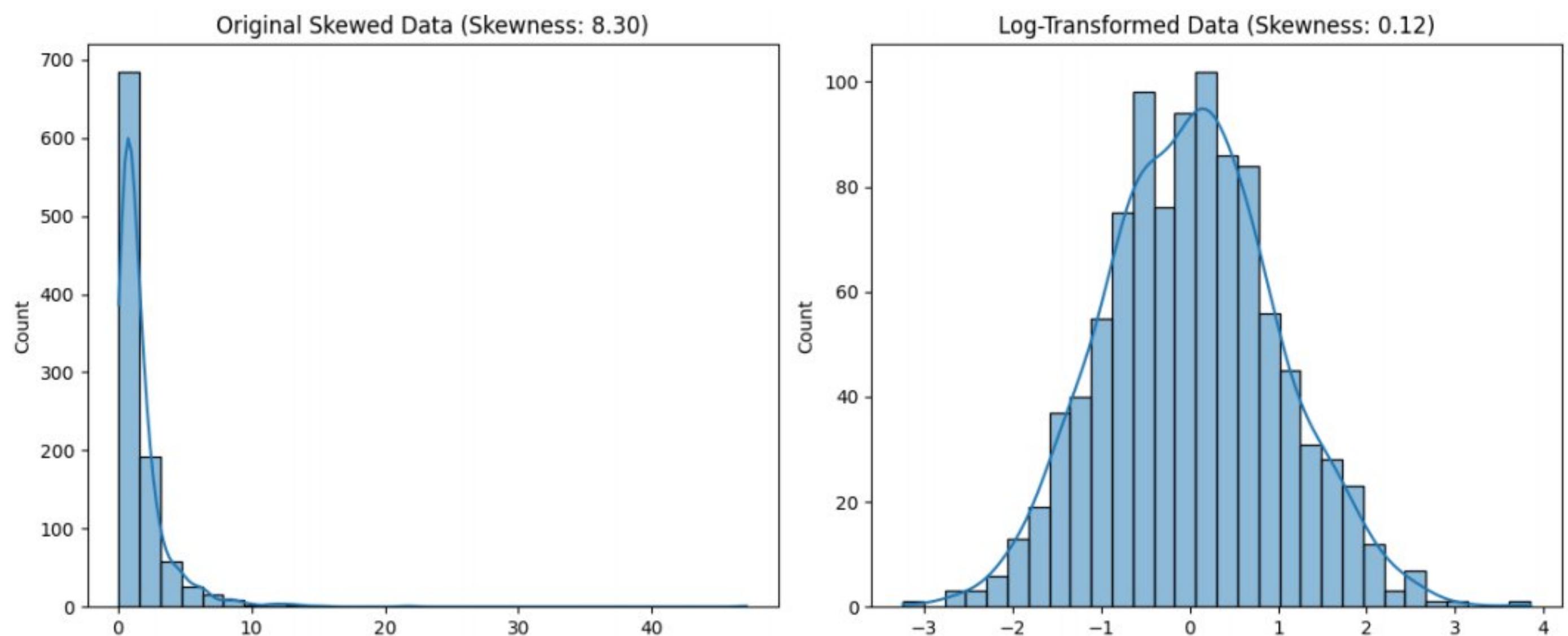
# Generate skewed data (log-normal)
np.random.seed(42)
skewed_data = np.random.lognormal(mean=0, sigma=1, size=1000)

# Original skewed data
plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
sns.histplot(skewed_data, kde=True, bins=30)
plt.title('Original Skewed Data (Skewness: {:.2f})'.format(stats.skew(skewed_d

# Apply log transformation to remove skew
log_transformed = np.log(skewed_data + 1e-8) # Add epsilon to avoid log(0)

plt.subplot(1, 2, 2)
sns.histplot(log_transformed, kde=True, bins=30)
plt.title('Log-Transformed Data (Skewness: {:.2f})'.format(stats.skew(log_tran

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



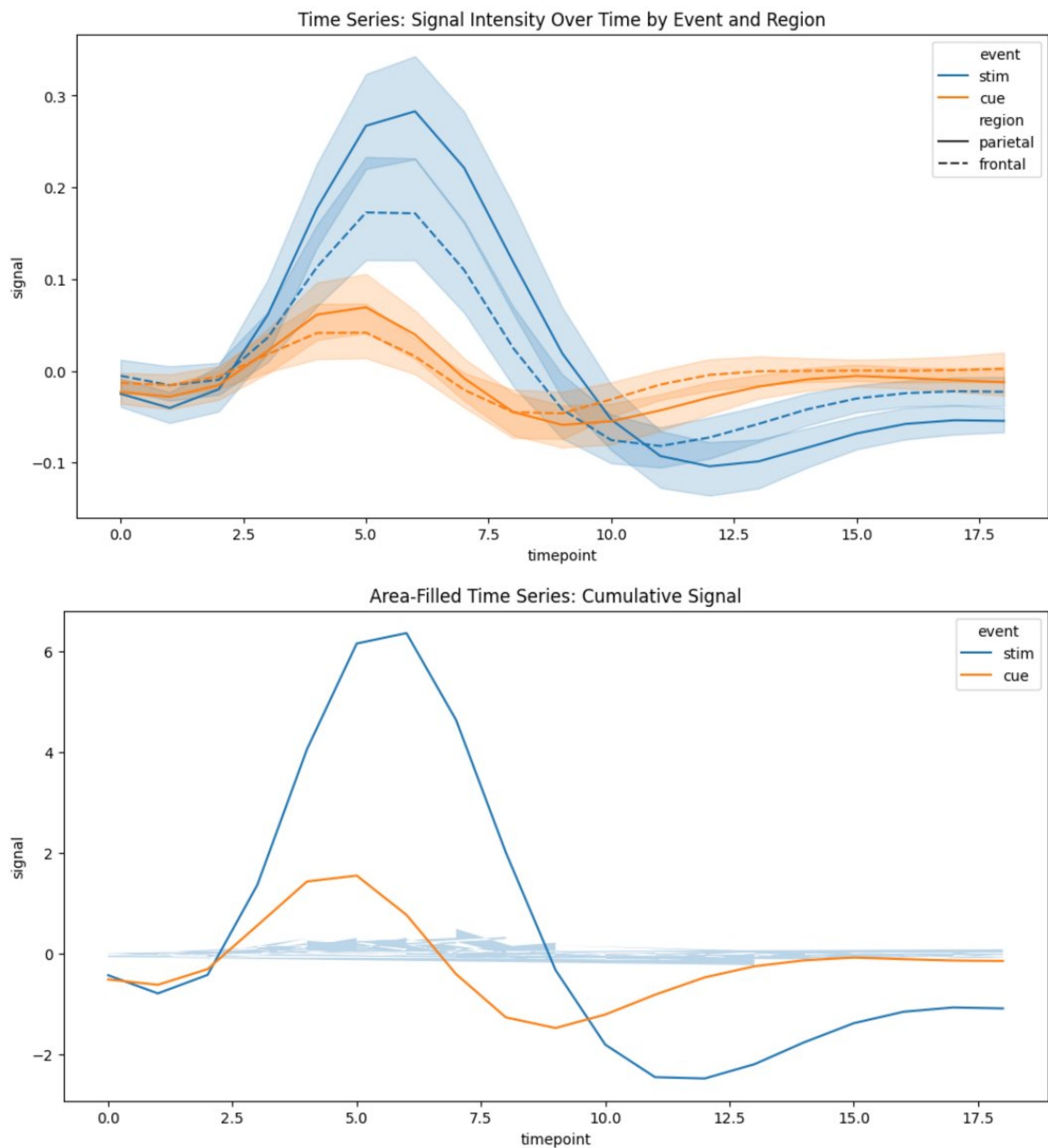
6. For a Sales Dataset, Do Time Series Visualization

```
In [ ]: import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd

# Load dataset (proxy: fmri time series)
fmri = sns.load_dataset('fmri')

# For real sales: df = pd.read_csv('SampleSuperstore.csv'); df['Order Date'] =
# Time Series Line Plot: Signal over time by event
plt.figure(figsize=(12, 6))
sns.lineplot(data=fmri, x='timepoint', y='signal', hue='event', style='region')
plt.title('Time Series: Signal Intensity Over Time by Event and Region')
plt.show()

# Area Plot for cumulative
plt.figure(figsize=(12, 6))
sns.lineplot(data=fmri, x='timepoint', y='signal', hue='event', estimator='sum')
plt.fill_between(fmri['timepoint'], fmri['signal'], alpha=0.3)
plt.title('Area-Filled Time Series: Cumulative Signal')
plt.show()
```



7. Build a Scatterplot and Suggest Dimension Reduction

```
In [ ]: import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
import pandas as pd

# Load dataset
iris = sns.load_dataset('iris')

# Scatterplot Matrix (pairplot for multi-scatter)
```

```
sns.pairplot(iris, hue='species')
plt.suptitle('Scatterplot Matrix: Iris Dataset', y=1.02)
plt.show()

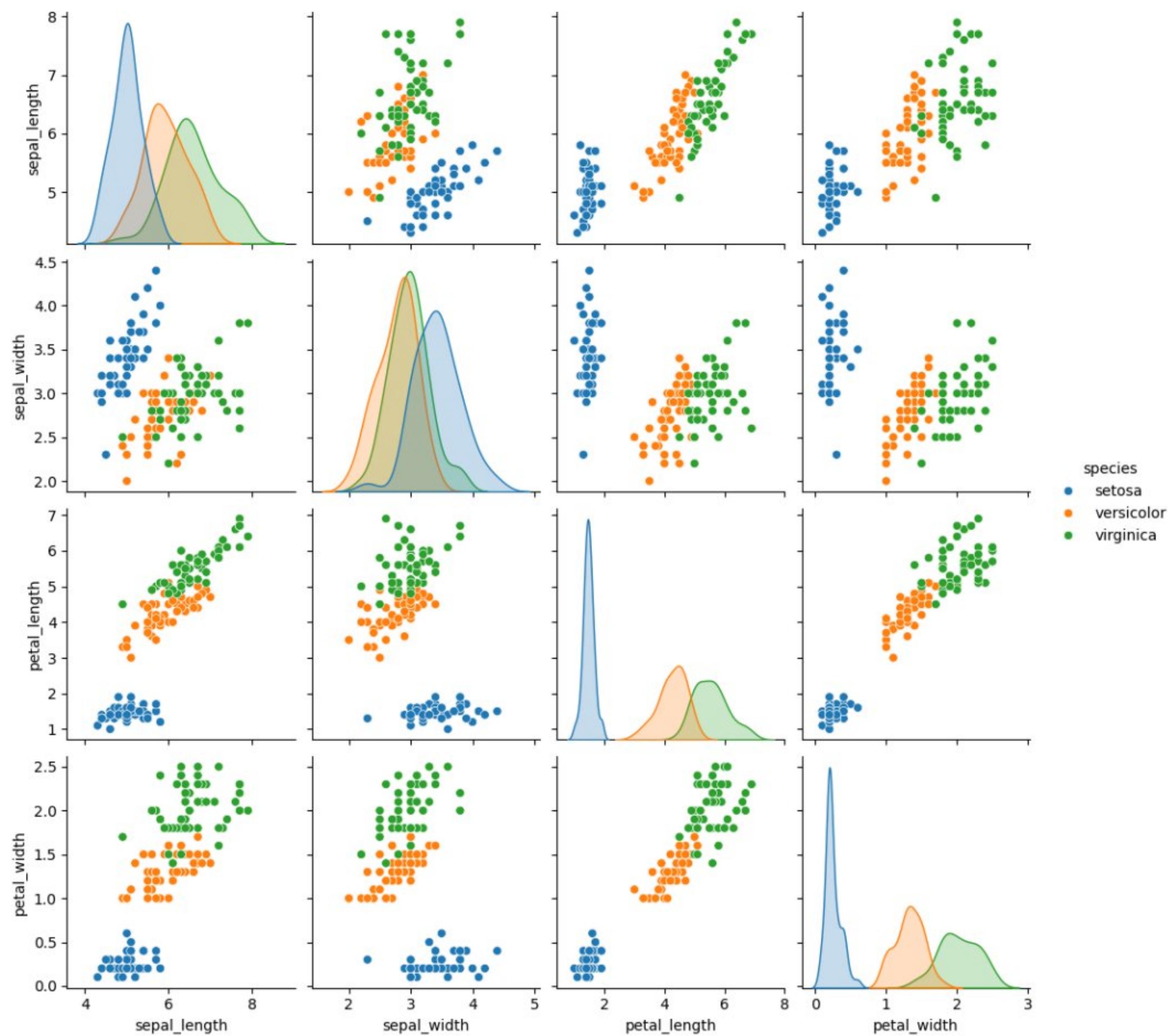
# Suggest Dimension Reduction: PCA to 2D
features = iris.select_dtypes(include='number').dropna()
pca = PCA(n_components=2)
pca_result = pca.fit_transform(features)

# Add to dataframe
pca_df = pd.DataFrame(pca_result, columns=['PC1', 'PC2'])
pca_df['species'] = iris['species'].values

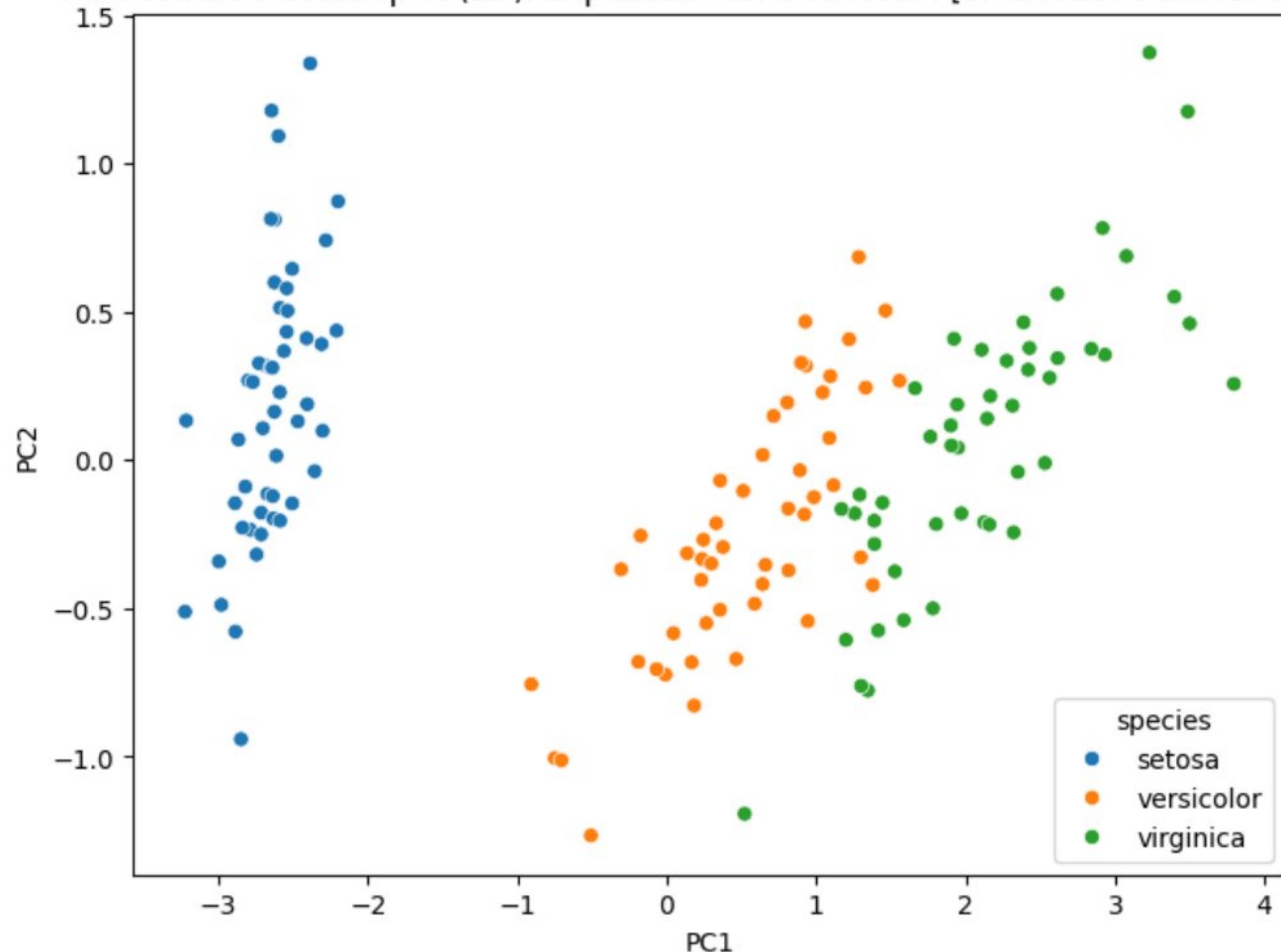
# Reduced Scatterplot
plt.figure(figsize=(8, 6))
sns.scatterplot(data=pca_df, x='PC1', y='PC2', hue='species')
plt.title('PCA-Reduced Scatterplot (2D): Explained Variance Ratio: {}'.format(
plt.show()

# Suggestion: Use PCA for dimensionality reduction when features > 2-3, to vis
```

Scatterplot Matrix: Iris Dataset



PCA-Reduced Scatterplot (2D): Explained Variance Ratio: [0.92461872 0.05306648]



In []: