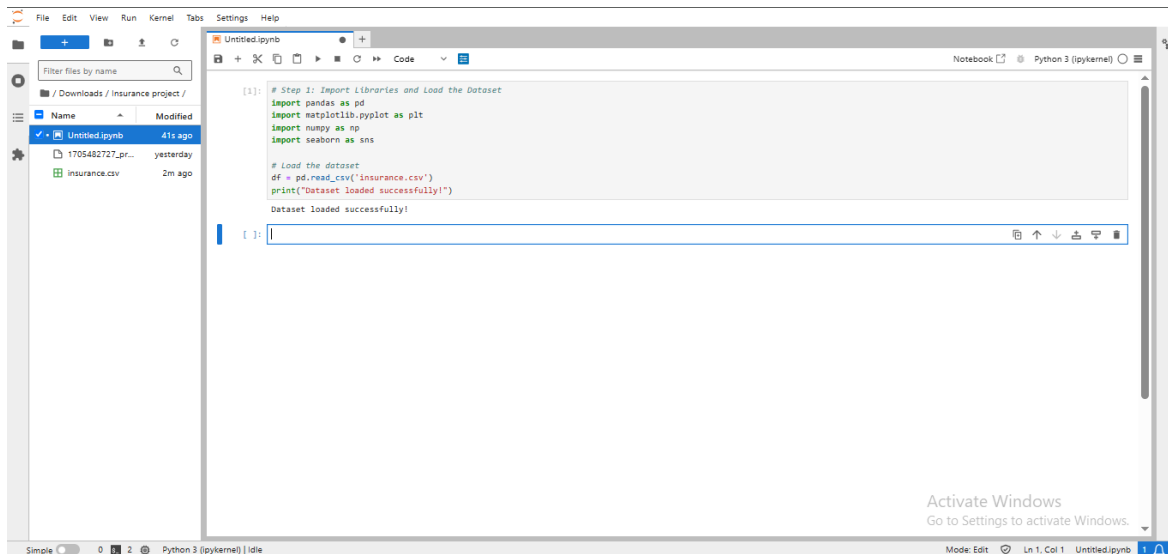**Insurance Data Analysis: Key Observations**

**Step 1: Import Libraries and Load the Dataset**
- Successfully imported essential libraries (Pandas, Matplotlib, NumPy, Seaborn).
- The insurance.csv dataset was loaded into a Pandas DataFrame without errors, confirming successful initial setup.



**Step 2: Check the Shape of the Data and Data Types**
- The dataset contains **1338** rows (representing policyholders) and **7** columns (representing features).
- Data types for each column are appropriate:
    - age (int)
    - sex (object)
    - bmi (float)
    - children (int)
    - smoker (object)
    - region (object)
    - charges (float)
- All columns have an equal number of non-null entries, indicating no immediate missing values.

## Step 3: Check Missing Values in the Dataset

- No missing values were found in any of the columns (df.isnull().sum() showed all zeros). The dataset is clean in terms of missing data, and no imputation steps are required.



## Step 4: Explore the Relationship Between Features and Target Column (charges)
## A. Count Plots for Categorical Features (Distribution)

- **Sex:** The distribution between male and female policyholders is relatively balanced.
- **Smoker Status:** A significant majority of policyholders are non-smokers, with a smaller proportion being smokers.
- **Region:** Policyholders are quite evenly distributed across the four geographical regions (northeast, northwest, southeast, southwest).

**B. Box Plots for Categorical Features vs. Charges**

- **Smoker:** There is a stark and significant difference in insurance charges between smokers and non-smokers. Smokers consistently face much higher premiums, indicating that smoking status is a primary determinant of insurance cost.
- **Sex:** The median charges for males and females are relatively similar, with no dramatic difference based solely on gender.
- **Region:** While minor variations exist, the median charges across different regions are quite similar, with no single region showing an extremely high or low premium compared to others.



**C. Scatter Plots for Numerical Features vs. Charges (Colored by Smoker Status)**

- **Age vs. Charges:** There's a clear positive linear relationship. As age increases, insurance premiums generally increase for both smokers and non-smokers. Two distinct bands of data points are visible, separating smokers (higher charges) from non-smokers (lower charges).
- **BMI vs. Charges:** A general positive trend is observed where higher BMI tends to be associated with higher charges, especially noticeable within both the smoker and non-smoker groups. However, the relationship is more scattered than with age.
- **Children vs. Charges:** The number of children does not show a strong direct linear correlation with charges. The data points are widely dispersed, suggesting that this feature alone is not a major predictor of premium costs.

Untitled.ipynb

Notebook ⬀   Python 3 (ipykernel) ◯ ☰

```python
[12]:   # Step 4C: Scatter plots of numerical columns vs. Charges (Fixed 'BMI' to 'bmi')
        plt.figure(figsize=(18, 6)) # Adjust figure size for better visualization

        plt.subplot(1, 3, 1) # 1 row, 3 columns, 1st plot
        sns.scatterplot(x='age', y='charges', data=df, hue='smoker', palette='coolwarm', alpha=0.7, s=50)
        plt.title('Charges vs. Age (Colored by Smoker Status)')
        plt.xlabel('Age')
        plt.ylabel('Charges')

        plt.subplot(1, 3, 2) # 1 row, 3 columns, 2nd plot
        # --- FIX APPLIED HERE: Changed x='BMI' to x='bmi' ---
        sns.scatterplot(x='bmi', y='charges', data=df, hue='smoker', palette='coolwarm', alpha=0.7, s=50)
        plt.title('Charges vs. BMI (Colored by Smoker Status)')
        plt.xlabel('BMI') # Label can remain 'BMI' for readability in the plot
        plt.ylabel('Charges')

        plt.subplot(1, 3, 3) # 1 row, 3 columns, 3rd plot
        sns.scatterplot(x='children', y='charges', data=df, hue='smoker', palette='coolwarm', alpha=0.7, s=50)
        plt.title('Charges vs. Number of Children (Colored by Smoker Status)')
        plt.xlabel('Number of Children')
        plt.ylabel('Charges')

        plt.tight_layout() # Adjust Layout to prevent overlapping titles/labels
        plt.show()
```

---

Untitled.ipynb

Notebook ⬀   Python 3 (ipykernel) ◯ ☰

```python
        plt.title('Charges vs. BMI (Colored by Smoker Status)')
        plt.xlabel('BMI') # Label can remain 'BMI' for readability in the plot
        plt.ylabel('Charges')

        plt.subplot(1, 3, 3) # 1 row, 3 columns, 3rd plot
        sns.scatterplot(x='children', y='charges', data=df, hue='smoker', palette='coolwarm', alpha=0.7, s=50)
        plt.title('Charges vs. Number of Children (Colored by Smoker Status)')
        plt.xlabel('Number of Children')
        plt.ylabel('Charges')

        plt.tight_layout() # Adjust Layout to prevent overlapping titles/labels
        plt.show()
```



```
[ ]:
```

Click to add a cell.

**Step 5: Perform Data Visualization Using Plots of Feature vs. Feature**

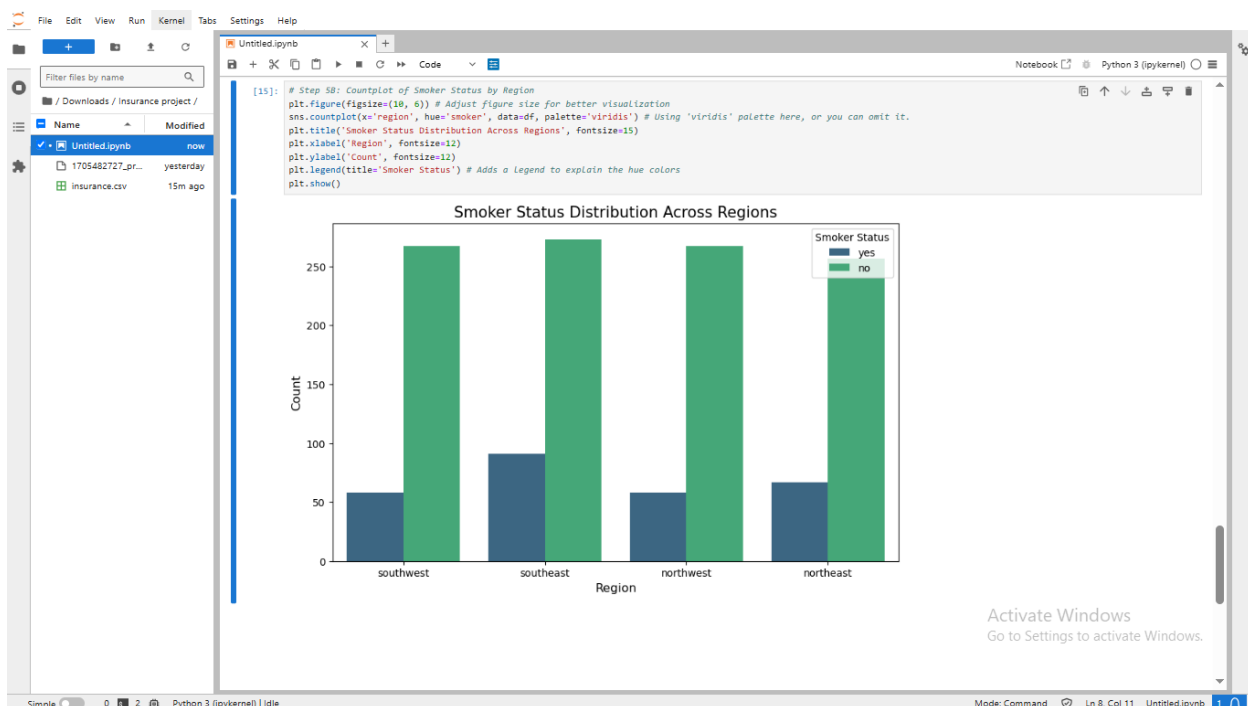**A. Correlation Matrix of Numerical Features**

- The correlation heatmap shows moderate positive correlations between charges and age (~0.30) and bmi (~0.20), confirming observations from scatter plots.
- age, bmi, and children show very low linear correlation with each other, indicating no significant multicollinearity among these independent numerical features.
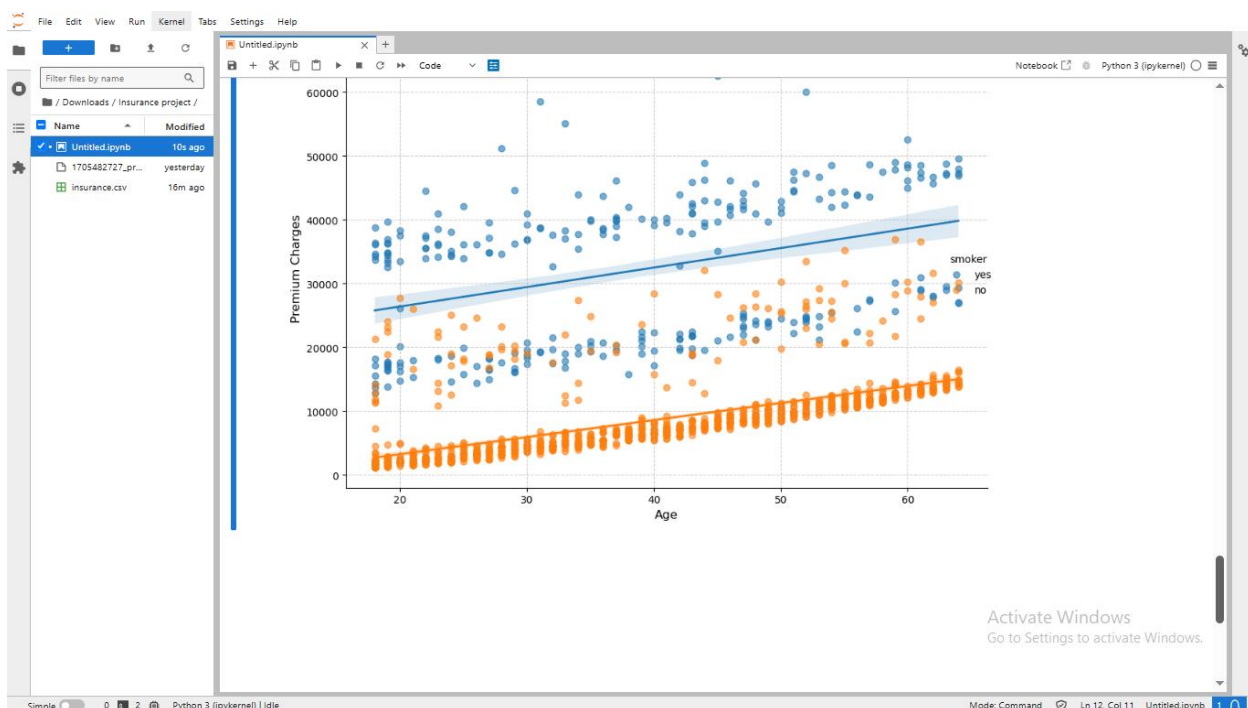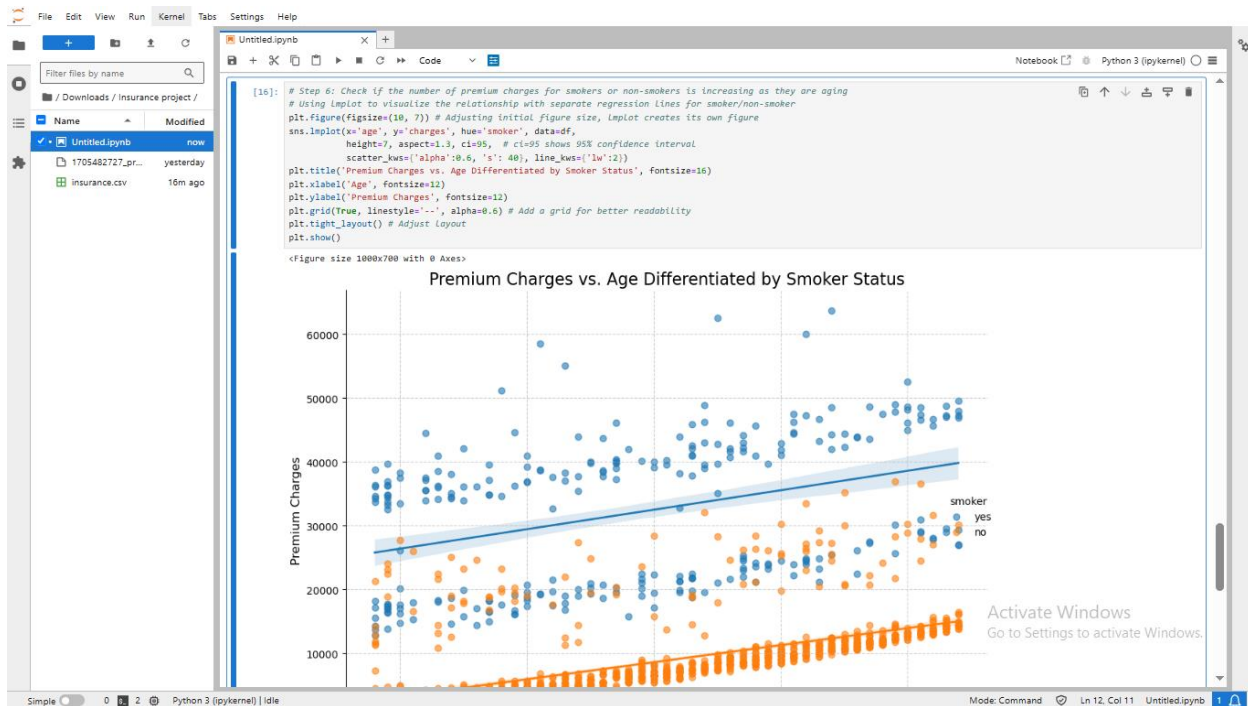


**B. Categorical Feature Relationships (Smoker Status by Region)**

- The proportion of smokers to non-smokers is relatively consistent across all four geographical regions, meaning no particular region has a disproportionately higher or lower number of smokers.

**Step 6: Check if the Number of Premium Charges for Smokers or Non-Smokers is Increasing as they are Aging**

- The Implot clearly demonstrates two distinct groups of charges based on smoker status.
- **Smokers:** Individuals who smoke consistently face significantly higher premium charges across all age groups compared to non-smokers.
- **Age Trend:** Within both the smoker and non-smoker groups, there is a clear positive linear trend: as age increases, the premium charges also increase.
- **Slope Comparison:** The regression line for smokers appears to have a slightly steeper slope than that for non-smokers, suggesting that the increase in charges with age might be somewhat more pronounced for smokers.

**Conclusion (Summary of Key Findings)**

This exploratory data analysis of the insurance dataset has provided valuable insights into the factors influencing medical insurance premiums. My/Our analysis revealed several key determinants:

- **Smoker Status** is the most significant predictor of insurance charges. I/We observed that smokers consistently incur substantially higher premiums compared to non-smokers, a difference that holds true across all age groups. This finding strongly supports the agency's initial hypothesis regarding the impact of smoking on chronic disease risk and associated healthcare costs.
- **Age** exhibits a clear positive linear relationship with insurance charges. As individuals age, their premiums tend to increase, a trend I/we found evident in both smoker and non-smoker groups, albeit with potentially a slightly steeper increase for smokers.
- **BMI** shows a positive correlation with charges, indicating that higher body mass index can lead to increased premiums, though this relationship is less linear and more dispersed than age.
- Features such as **sex, number of children**, and **geographical region** showed less direct or significant impact on the magnitude of insurance charges compared to smoker status and age. While minor variations exist, they are not as pronounced.

Overall, the dataset was found to be clean and complete, providing a robust foundation for further analysis. These insights are crucial for ABC Insurance to understand their current premium structure and identify key risk factors among their policyholders.

**Next Steps / Future Work**

The insights gained from this exploratory data analysis lay a strong foundation for building a predictive model for medical insurance charges. The logical next steps for this project would involve:

1. **Data Preprocessing:**
   - **Encoding Categorical Features:** Convert categorical variables such as 'sex', 'smoker', and 'region' into numerical representations (e.g., using One-Hot Encoding) so they can be processed by machine learning algorithms.
   - **Feature Scaling (Optional but Recommended):** Scale numerical features like 'age' and 'bmi' to ensure that no single feature dominates the model due to its larger numerical range.
2. **Model Building:** Select and train appropriate regression models (e.g., Linear Regression, Decision Trees, Random Forests, Gradient Boosting) using the prepared dataset.
3. **Model Evaluation:** Evaluate the performance of the trained model(s) using relevant metrics (e.g., R-squared, Mean Absolute Error, Root Mean Squared Error) to determine its accuracy and predictive power in forecasting medical insurance costs.

By following these steps, ABC Insurance can develop a robust model to predict insurance premiums, further enhancing their business decision-making and operational efficiency.