

AIR QUALITY TEAM E - EXTENSION

Zipcode - Analysis

```
In [44]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from tabulate import tabulate
from scipy import stats
import geopandas as gpd
from shapely import wkt
import plotly.graph_objects as go
from ipywidgets import widgets, Output
import pandas as pd
from ipywidgets import interact, Dropdown
```

```
In [45]: # Load the datasets
health_data_2022 = pd.read_csv('../Merged_Data/2022_Health_Zipcode.csv')
health_data_2023 = pd.read_csv('../Merged_Data/2023_Health_Zipcode.csv')
```

```
In [46]: health_data_2022
```

Out [46]:

	Zipcode	TotalPopulation	ACCESS2_CrudePrev	ARTHRITIS_CrudePrev	BINGE_Crud
0	2111	7383	7.8		16.1
1	2113	6915	4.2		11.3
2	2118	26498	8.9		16.6
3	2124	47783	10.8		21.9
4	2127	31799	6.4		15.7
5	2128	40508	16.3		17.5
6	2130	35401	7.2		17.3
7	2135	42780	6.1		13.4
8	2139	36349	5.5		12.1

9 rows × 35 columns

In [47]:

health_data_2023

Out[47]:

	Zipcode	TotalPopulation	ACCESS2_CrudePrev	ARTHRITIS_CrudePrev	BINGE_Crud
0	2111	7383	5.5		15.8
1	2113	6915	2.6		12.3
2	2118	26498	6.8		17.1
3	2124	47783	7.7		22.4
4	2127	31799	4.6		16.7
5	2128	40508	13.2		17.9
6	2130	35401	5.1		18.2
7	2135	42780	4.1		14.0
8	2139	36349	3.6		13.2

9 rows × 42 columns

In [48]:

```
pm25_2022 = health_data_2022[['Zipcode', 'PM2.5']]
pm25_2023 = health_data_2023[['Zipcode', 'PM2.5']]

pm25_comparison = pd.merge(pm25_2022,
                           pm25_2023,
                           on='Zipcode',
                           suffixes=('_2022', '_2023'))

pm25_comparison['PM2.5_Change'] = pm25_comparison[
    'PM2.5_2023'] - pm25_comparison['PM2.5_2022']

table = tabulate(pm25_comparison,
                 headers='keys',
                 tablefmt='fancy_grid',
                 showindex=False)
print(table)
```

Zipcode	PM2.5_2022	PM2.5_2023	PM2.5_Change
2111	6.73753	9.4645	2.72697
2113	7.9084	9.96361	2.05521
2118	8.30109	10.5777	2.27656
2124	7.21404	9.10869	1.89465
2127	7.75032	10.5325	2.78215
2128	7.08768	9.64689	2.55921
2130	7.12675	8.7106	1.58385
2135	6.94934	8.87072	1.92138
2139	5.53174	9.11105	3.57931

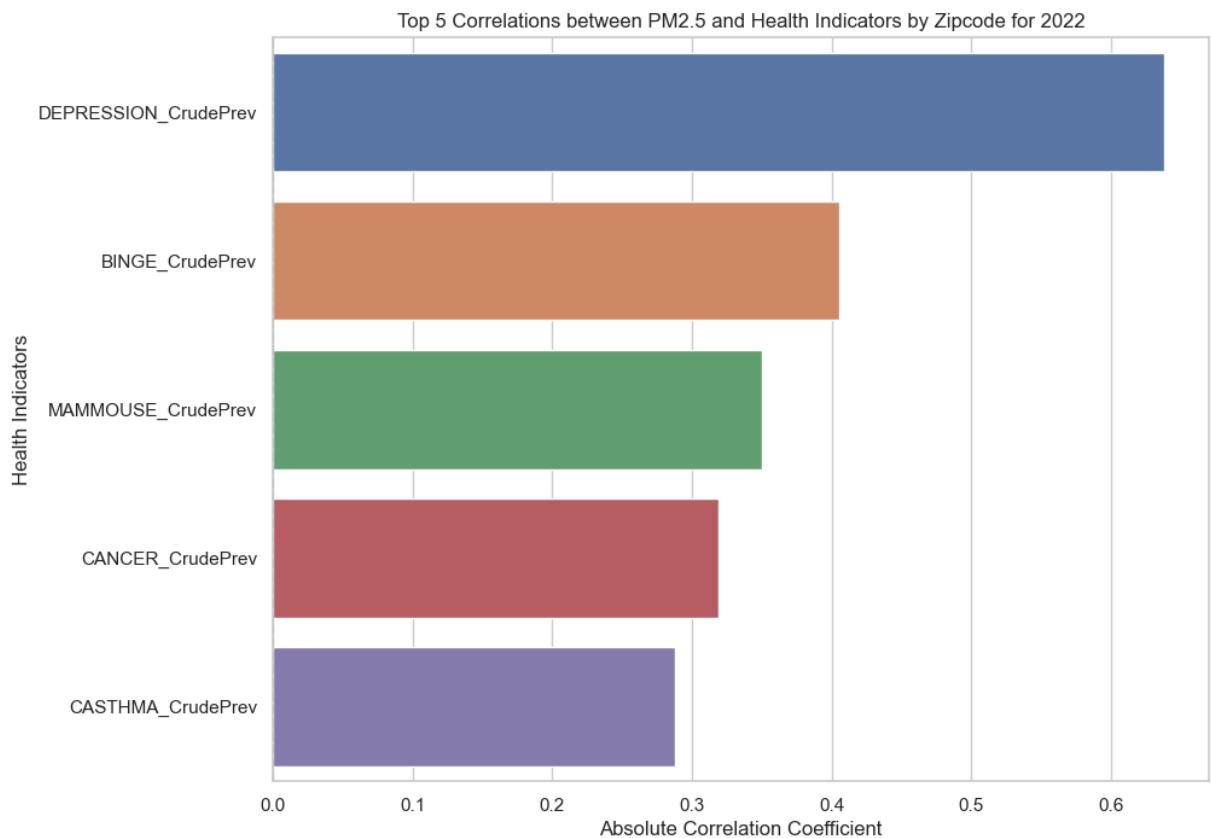
The comparison of air quality between 2022 and 2023, as indicated by PM2.5 levels across common zip codes, reveals a notable trend. For the zip codes analyzed, there is a consistent increase in PM2.5 levels from 2022 to 2023. This upward shift suggests a deterioration in air quality over the year. The changes vary across different areas, with some zip codes experiencing more significant increases than others.

```
In [53]: import matplotlib.pyplot as plt
import seaborn as sns
crude_columns_2022 = [
    col for col in health_data_2022.columns if 'Crude' in col
]
crude_columns_2022.append('PM2.5')
correlation_matrix_2022 = health_data_2022[crude_columns_2022].corr()

correlations_2022 = correlation_matrix_2022['PM2.5'].drop(
    'PM2.5') # excluding the PM2.5 self-correlation

# Getting the top 5 correlations by absolute value
top_5_correlations = correlations_2022.abs().nlargest(5)

# Creating a bar plot to visualize the top 5 correlations by zipcode
plt.figure(figsize=(10, 8))
sns.barplot(x=top_5_correlations.values, y=top_5_correlations.index)
plt.title('Top 5 Correlations between PM2.5 and Health Indicators by Zipcode')
plt.xlabel('Absolute Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```

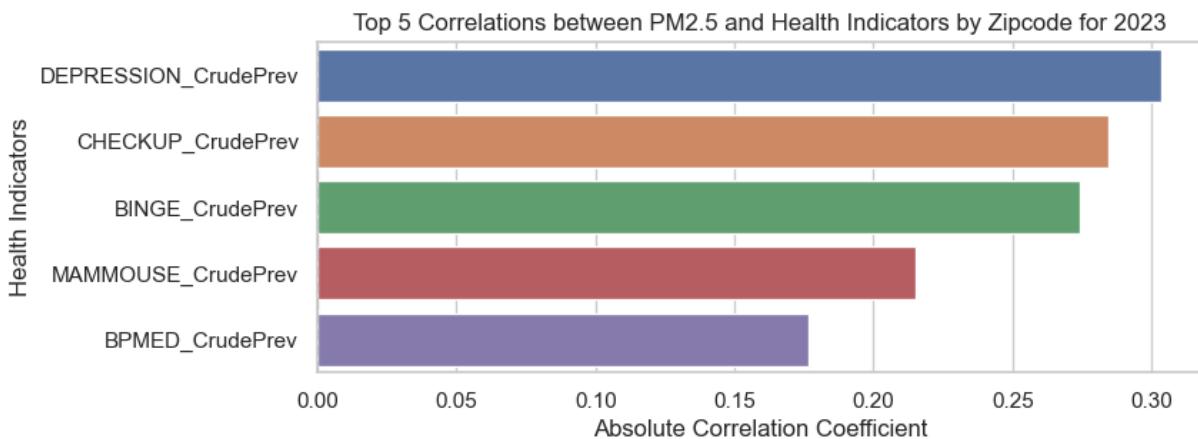


```
In [62]: crude_columns_2023 = [
    col for col in health_data_2023.columns if 'Crude' in col
]
crude_columns_2023.append('PM2.5')
correlation_matrix_2023 = health_data_2023[crude_columns_2023].corr()

correlations_2023 = correlation_matrix_2023['PM2.5'].drop(
    'PM2.5') # excluding the PM2.5 self-correlation

# Getting the top 5 correlations by absolute value
top_5_correlations = correlations_2023.abs().nlargest(5)

# Creating a bar plot to visualize the top 5 correlations by zipcode
plt.figure(figsize=(8, 3))
sns.barplot(x=top_5_correlations.values, y=top_5_correlations.index)
plt.title('Top 5 Correlations between PM2.5 and Health Indicators by Zipcode')
plt.xlabel('Absolute Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```



```
In [63]: import matplotlib.pyplot as plt
import seaborn as sns

# Assuming you have defined correlations_2022 and correlations_2023

# Getting the top 5 correlations by absolute value for 2022
top_5_correlations_2022 = correlations_2022.abs().nlargest(5)

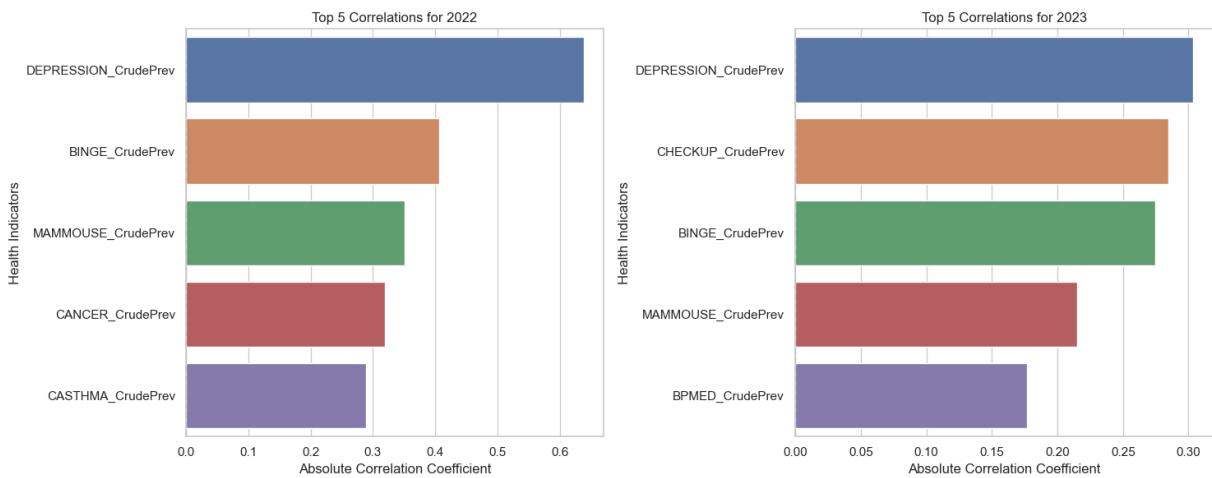
# Getting the top 5 correlations by absolute value for 2023
top_5_correlations_2023 = correlations_2023.abs().nlargest(5)

# Create subplots to display side-by-side bar plots for 2022 and 2023
fig, axes = plt.subplots(1, 2, figsize=(15, 6))

# Plotting top 5 correlations for 2022
sns.barplot(ax=axes[0], x=top_5_correlations_2022.values, y=top_5_correlations_2022.index)
axes[0].set_title('Top 5 Correlations for 2022')
axes[0].set_xlabel('Absolute Correlation Coefficient')
axes[0].set_ylabel('Health Indicators')
axes[0].axvline(x=0, color='gray', linestyle='--')

# Plotting top 5 correlations for 2023
sns.barplot(ax=axes[1], x=top_5_correlations_2023.values, y=top_5_correlations_2023.index)
axes[1].set_title('Top 5 Correlations for 2023')
axes[1].set_xlabel('Absolute Correlation Coefficient')
axes[1].set_ylabel('Health Indicators')
axes[1].axvline(x=0, color='gray', linestyle='--')

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



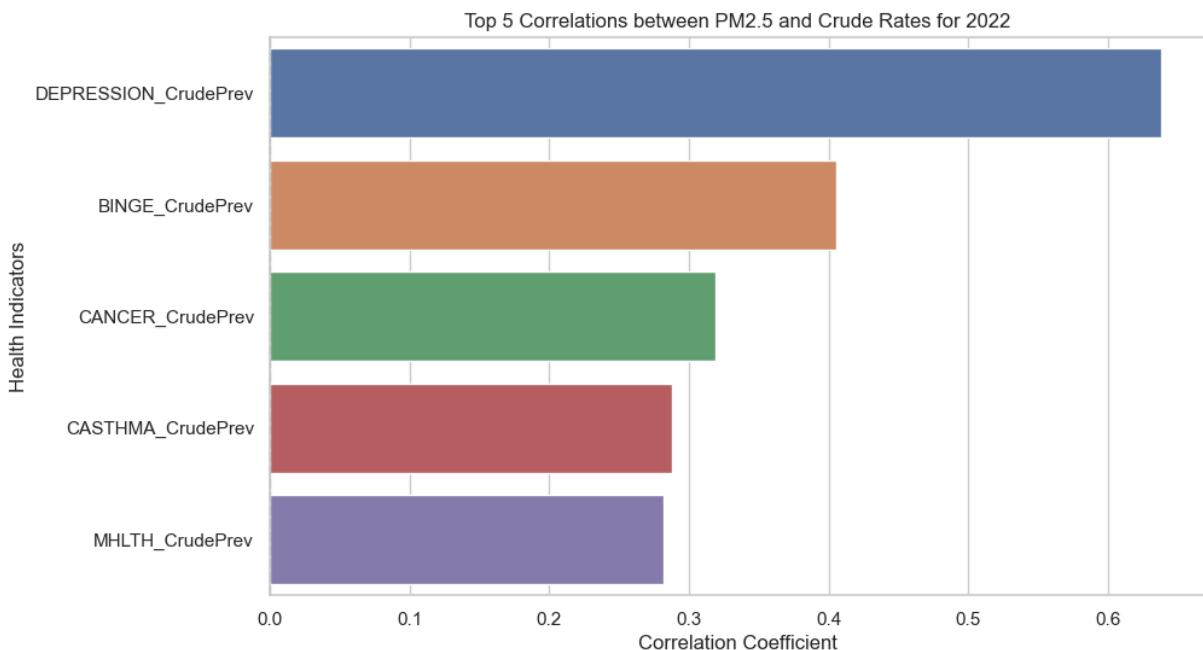
```
In [28]: import matplotlib.pyplot as plt
import seaborn as sns

crude_columns_2022 = [
    col for col in health_data_2022.columns if 'Crude' in col
]
crude_columns_2022.append('PM2.5')
correlation_matrix_2022 = health_data_2022[crude_columns_2022].corr()

correlations_2022 = correlation_matrix_2022['PM2.5'].drop('PM2.5')

# Getting the top 3 highest coefficients
top_3_correlations = correlations_2022.nlargest(5)

# Creating a bar plot to visualize the top 3 correlations
plt.figure(figsize=(10, 6))
sns.barplot(x=top_3_correlations.values, y=top_3_correlations.index)
plt.title('Top 5 Correlations between PM2.5 and Crude Rates for 2022')
plt.xlabel('Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```



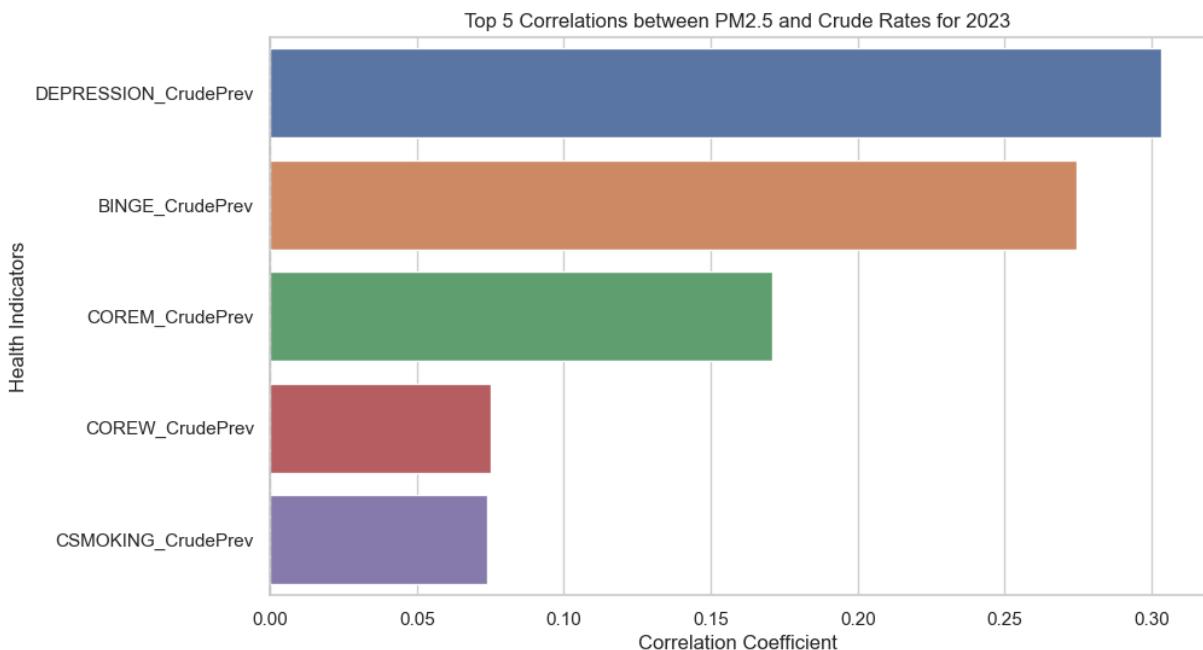
```
In [27]: import matplotlib.pyplot as plt
import seaborn as sns

crude_columns_2023 = [
    col for col in health_data_2023.columns if 'Crude' in col
]
crude_columns_2023.append('PM2.5')
correlation_matrix_2023 = health_data_2023[crude_columns_2023].corr()

correlations_2023 = correlation_matrix_2023['PM2.5'].drop('PM2.5')

# Getting the top 3 highest coefficients
top_3_correlations = correlations_2023.nlargest(5)

# Creating a bar plot to visualize the top 3 correlations
plt.figure(figsize=(10, 6))
sns.barplot(x=top_3_correlations.values, y=top_3_correlations.index)
plt.title('Top 5 Correlations between PM2.5 and Crude Rates for 2023')
plt.xlabel('Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```



```
In [36]: common_columns = set(health_data_2022.columns).intersection(set(health_data_2023.columns))

# Filtering out non-disease related columns
non_disease_columns = {'Zipcode', 'TotalPopulation', 'Geolocation', 'Year',
disease_columns = common_columns - non_disease_columns

# Merging the two datasets on Zipcode to calculate the changes
merged_data = pd.merge(health_data_2022[['Zipcode']] + list(disease_columns),
                      health_data_2023[['Zipcode']] + list(disease_columns),
                      on='Zipcode', suffixes=('_2022', '_2023'))

# Calculating the change for each disease and PM2.5
for disease in disease_columns:
    merged_data[f'{disease}_Change'] = merged_data[f'{disease}_2023'] - merged_data[f'{disease}_2022']

# Calculating PM2.5 change
merged_data['PM2.5_Change'] = merged_data['PM2.5_2023'] - merged_data['PM2.5_2022']

# Selecting only the change columns along with the Zipcode for the final table
change_columns = [col for col in merged_data.columns if 'Change' in col]
final_table = merged_data[['Zipcode']] + change_columns

final_table
```

Out[36]:

	Zipcode	DIABETES_CrudePrev_Change	BINGE_CrudePrev_Change	STROKE_CrudeP
0	2111	0.9	0.1	
1	2113	0.1	0.1	
2	2118	0.5	-0.2	
3	2124	0.3	0.2	
4	2127	0.2	0.0	
5	2128	0.2	-0.1	
6	2130	0.2	0.2	
7	2135	0.3	-0.1	
8	2139	0.1	2.0	

9 rows × 33 columns

In [38]:

```
# Assuming the 'final_table' DataFrame structure looks like the provided data

selected_columns = ['DIABETES_CrudePrev_Change', 'BINGE_CrudePrev_Change', 'COREM_CrudePrev_Change', 'CHECKUP_CrudePrev_Change', 'COPD_CrudePrev_Change', 'CSMOKING_CrudePrev_Change', 'CHD_CrudePrev_Change', 'CHF_CrudePrev_Change', 'HBP_CrudePrev_Change', 'PM2.5_Change']

# Add other columns related to the top 5 diseases here
#selected_columns.append('PM2.5_Change')

final_table = final_table[selected_columns]
final_table
```

Out[38]:

	DIABETES_CrudePrev_Change	BINGE_CrudePrev_Change	STROKE_CrudePrev_Change	COPD_CrudePrev_Change
0	0.9	0.1	0.1	0.0
1	0.1	0.1	0.1	0.0
2	0.5	-0.2	-0.2	0.0
3	0.3	0.2	0.2	0.0
4	0.2	0.0	0.0	0.0
5	0.2	-0.1	-0.1	0.0
6	0.2	0.2	0.2	0.0
7	0.3	-0.1	-0.1	0.0
8	0.1	2.0	0.0	0.0

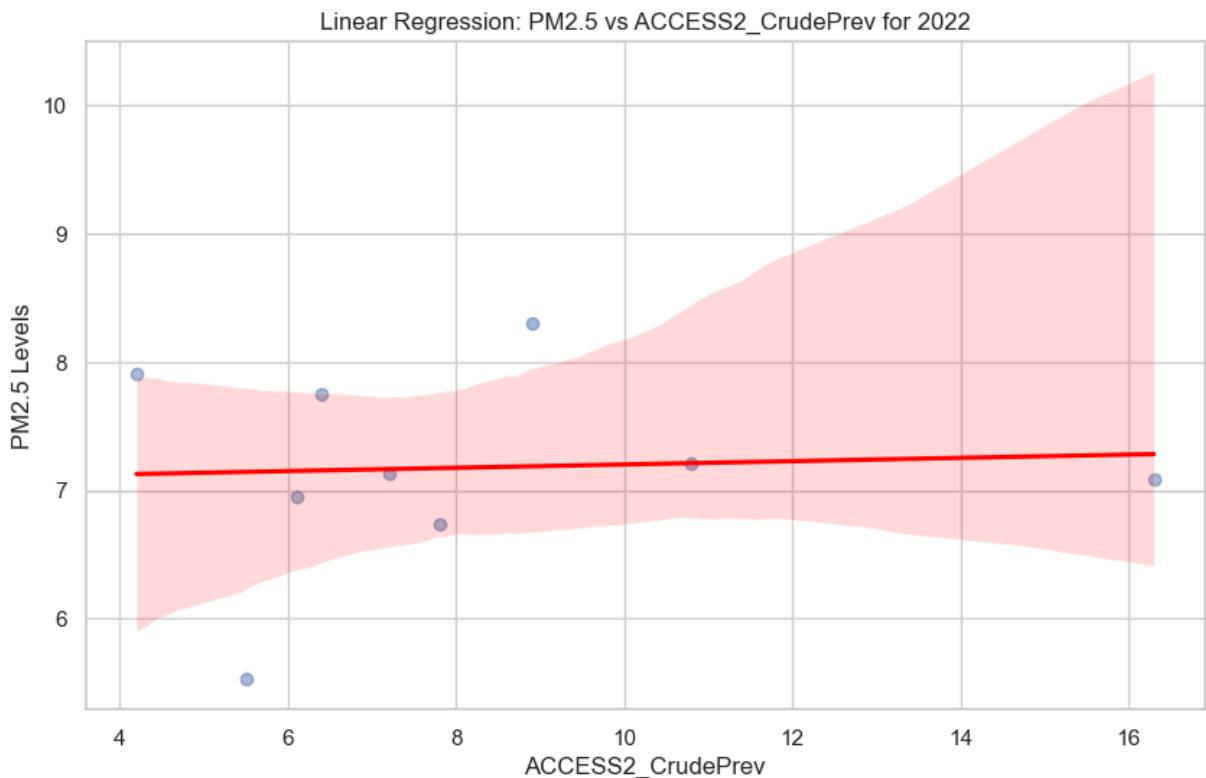
In [39]:

```
def plot_linear_regression(data, crude_columns, year):
    for column in crude_columns:
        if column != 'PM2.5':
            # Performing linear regression
            slope, intercept, r_value, p_value, std_err = stats.linregress(data['Year'], data[column])
            print(f"Linear regression for {column}: Slope = {slope}, Intercept = {intercept}, R-value = {r_value}, P-value = {p_value}, Std Err = {std_err}")
```

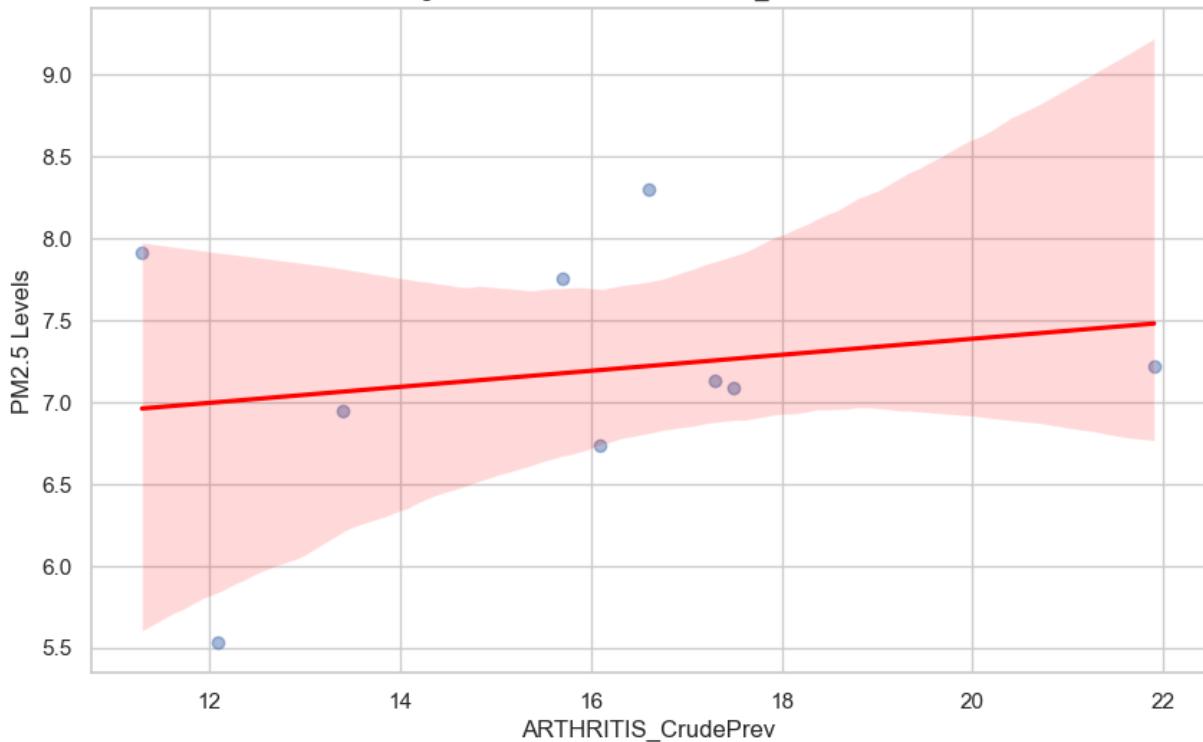
```
# Plotting
plt.figure(figsize=(10, 6))
sns.regplot(x=column, y='PM2.5', data=data, scatter_kws={'alpha': 0.5})
plt.title(f'Linear Regression: PM2.5 vs {column} for {year}')
plt.xlabel(column)
plt.ylabel('PM2.5 Levels')
plt.show()

# Extracting crude columns from 2022 dataset
crude_columns_2022 = [col for col in health_data_2022.columns if 'Crude' in col]

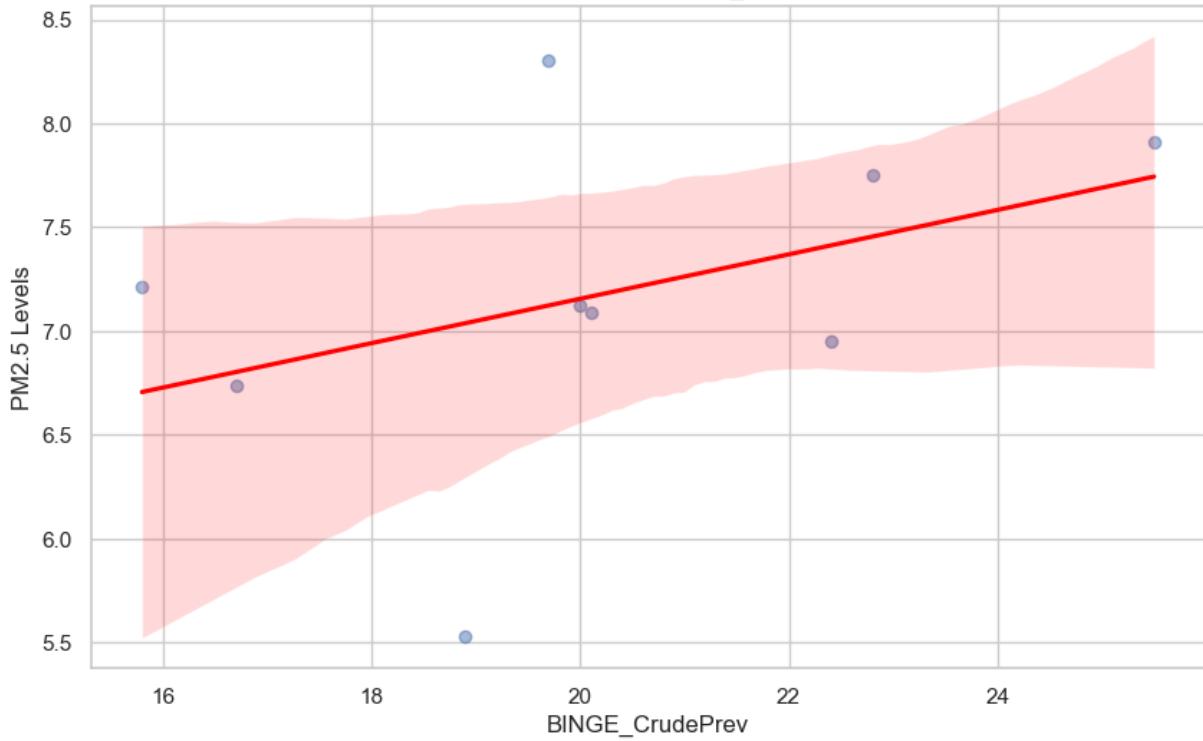
# Making linear regression plots for each crude rate against PM2.5 level for 2022
plot_linear_regression(health_data_2022, crude_columns_2022, '2022')
```



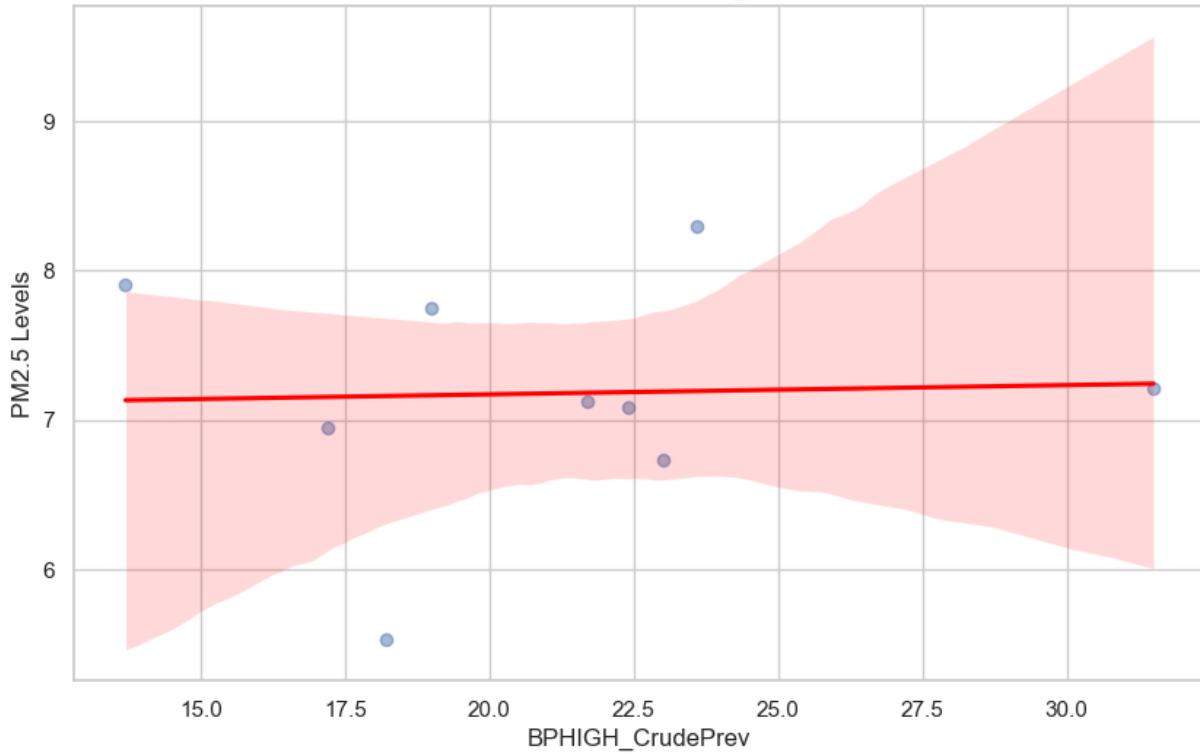
Linear Regression: PM2.5 vs ARTHRITIS_CrudePrev for 2022



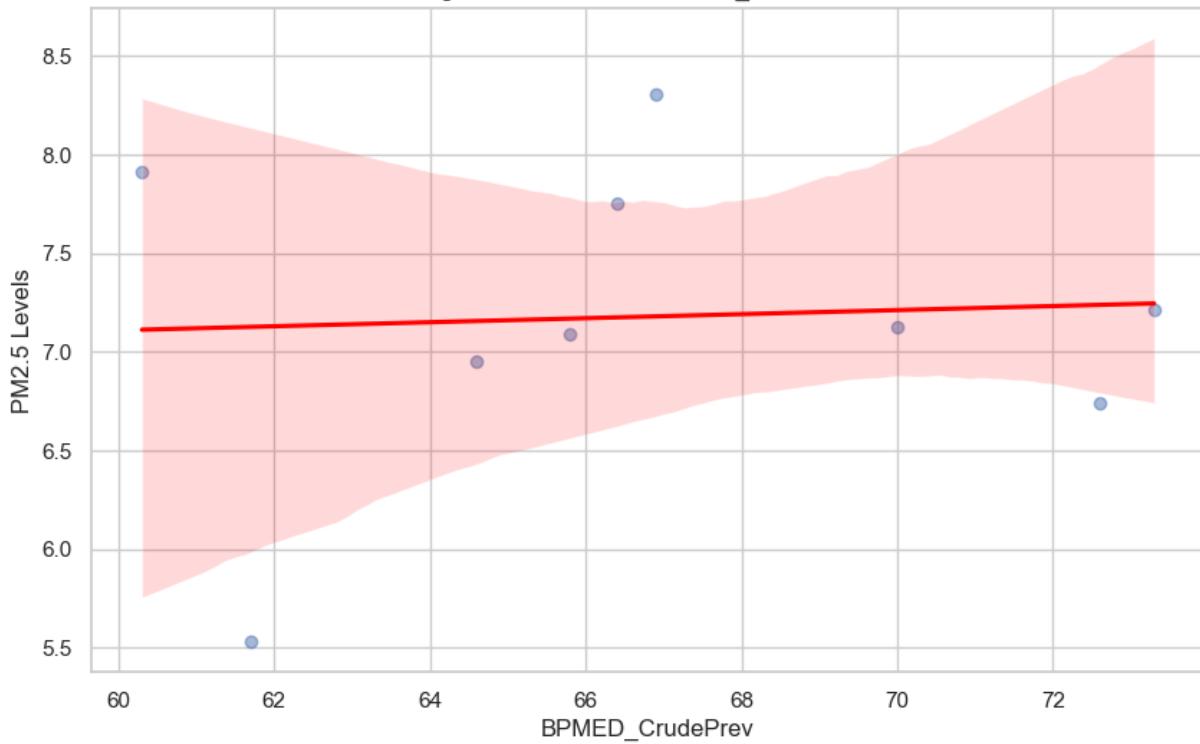
Linear Regression: PM2.5 vs BINGE_CrudePrev for 2022



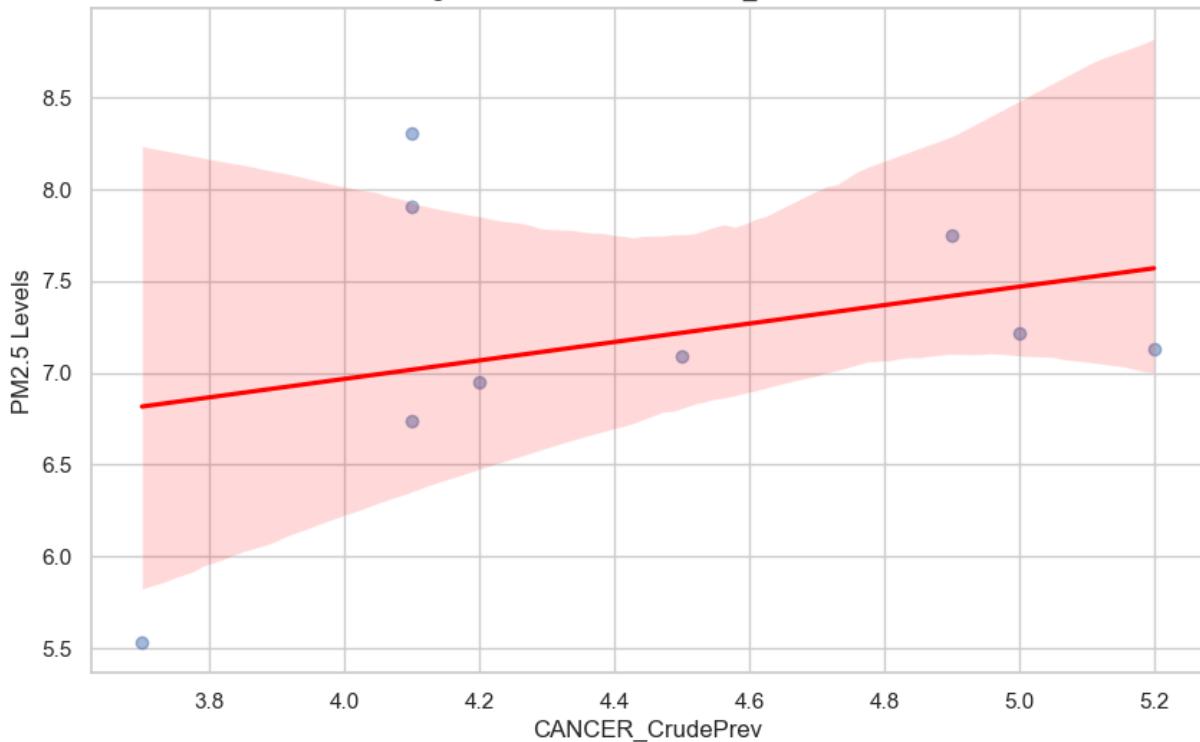
Linear Regression: PM2.5 vs BPHIGH_CrudePrev for 2022



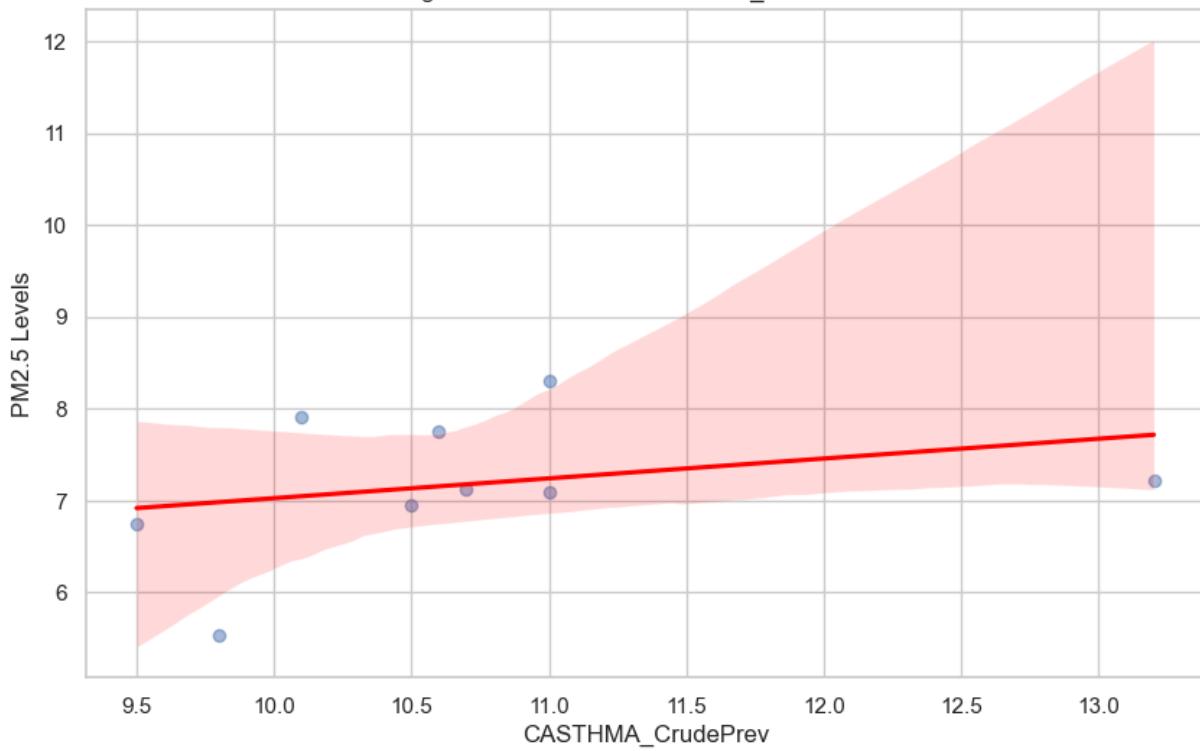
Linear Regression: PM2.5 vs BPMED_CrudePrev for 2022



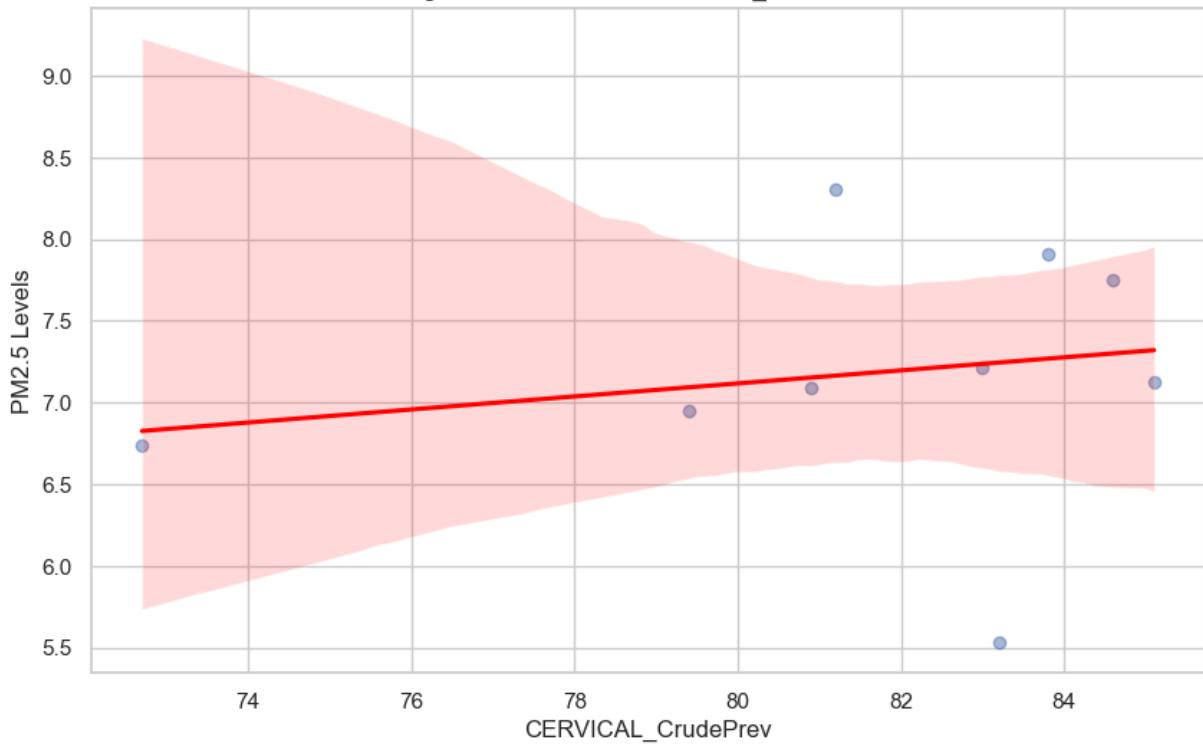
Linear Regression: PM2.5 vs CANCER_CrudePrev for 2022



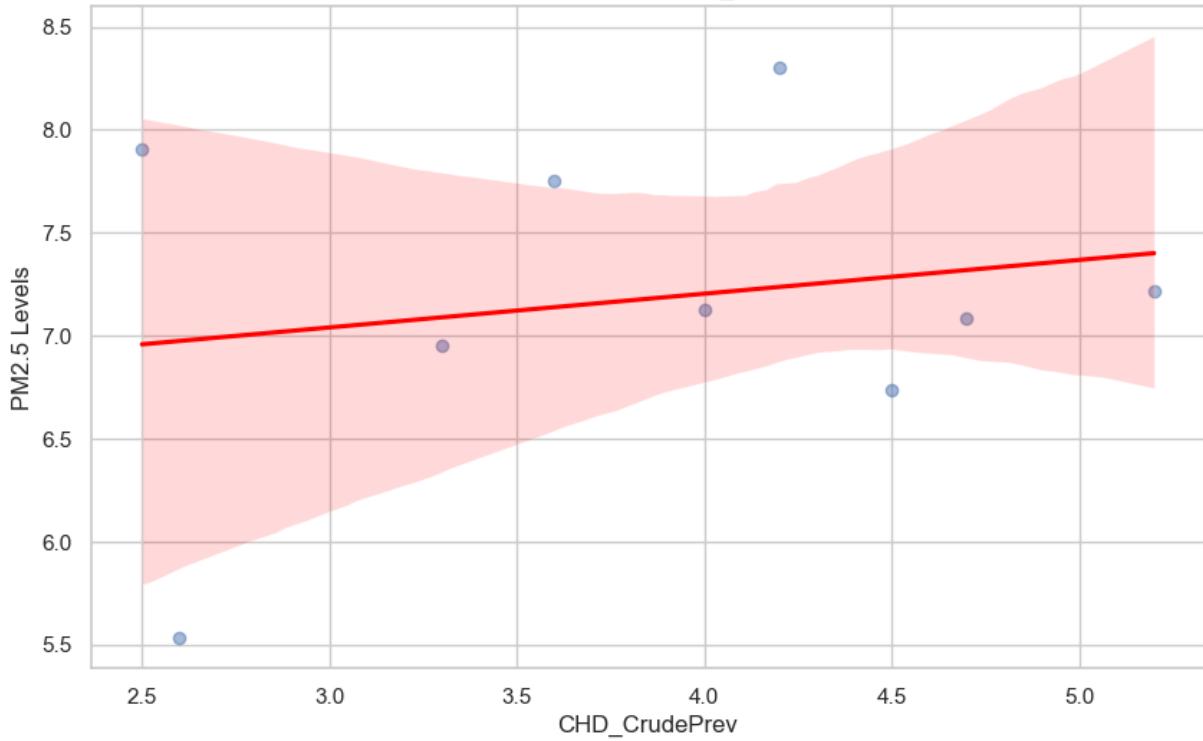
Linear Regression: PM2.5 vs CASTHMA_CrudePrev for 2022



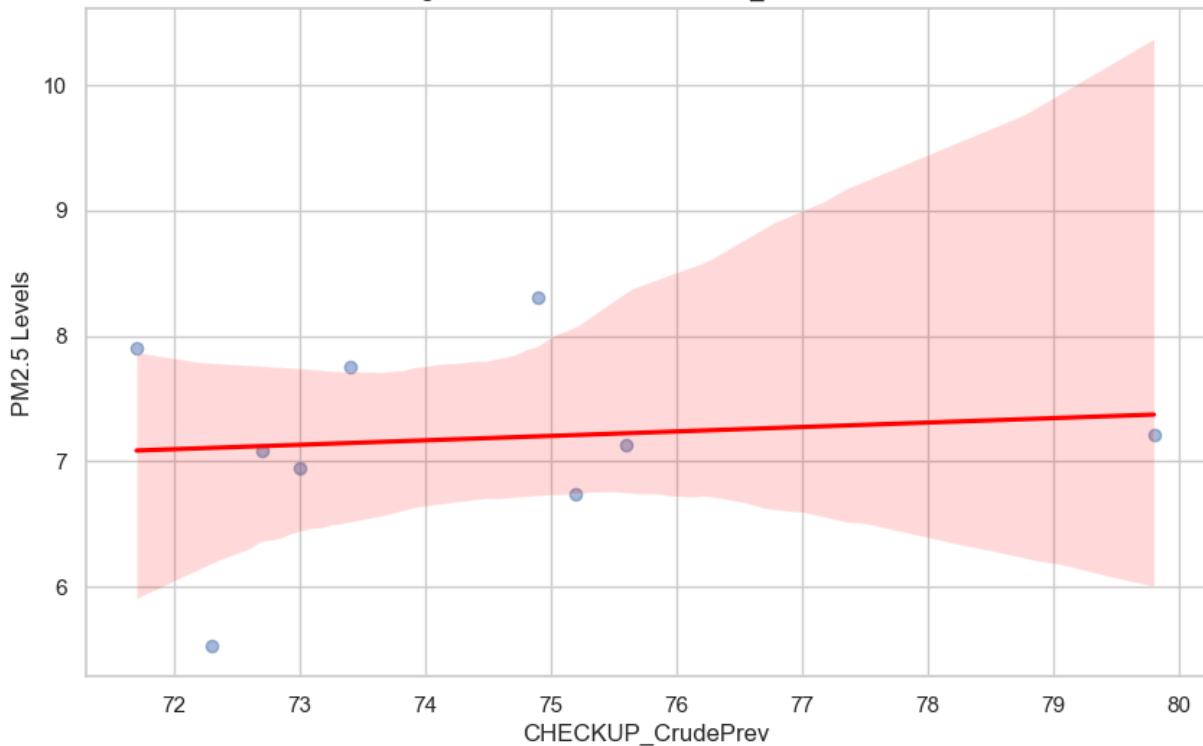
Linear Regression: PM2.5 vs CERVICAL_CrudePrev for 2022



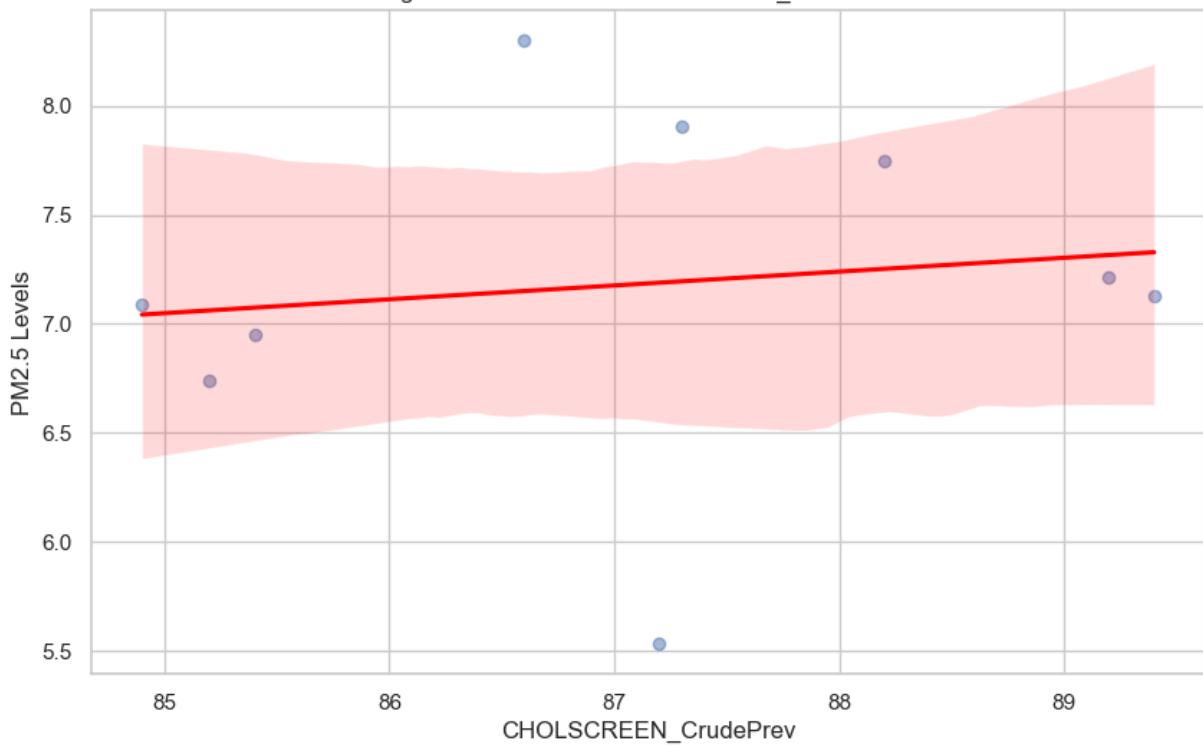
Linear Regression: PM2.5 vs CHD_CrudePrev for 2022



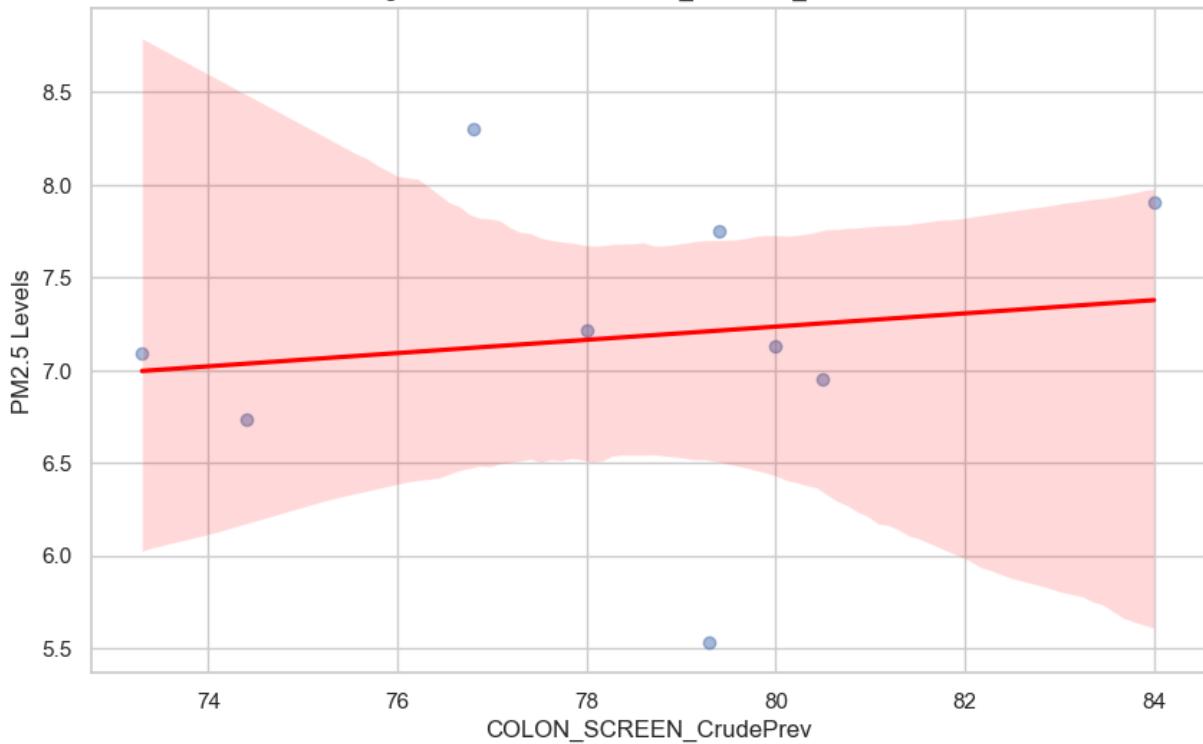
Linear Regression: PM2.5 vs CHECKUP_CrudePrev for 2022



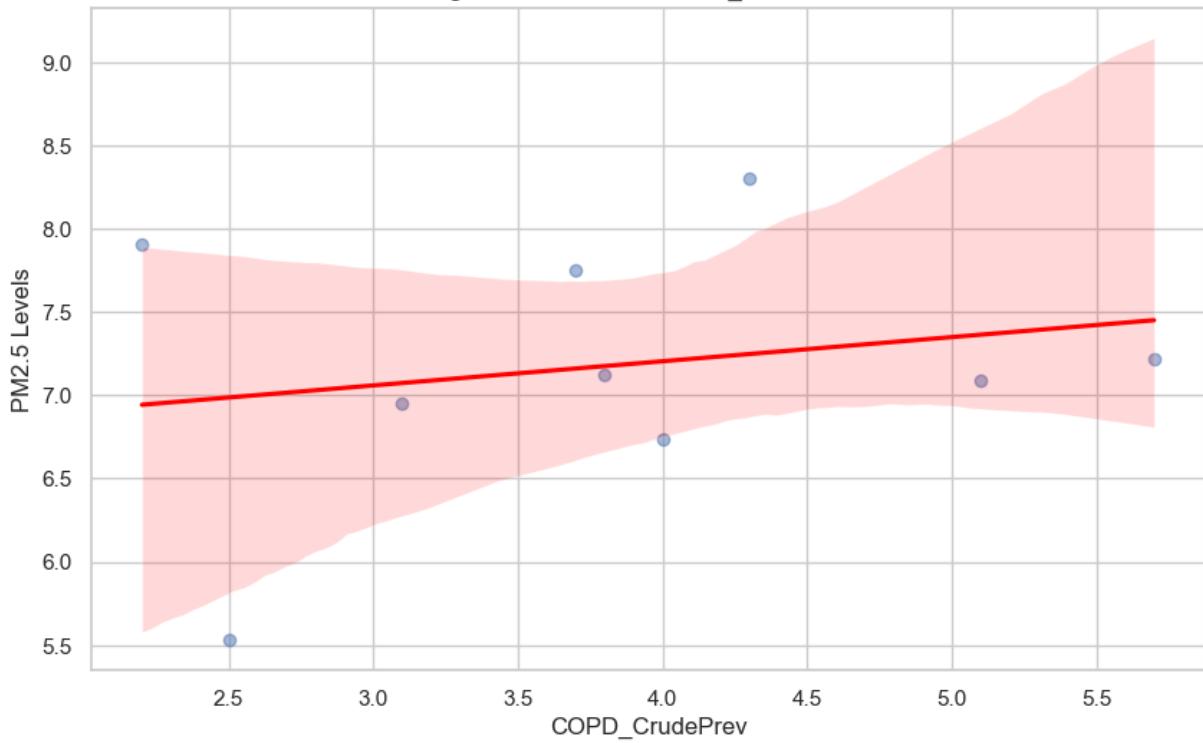
Linear Regression: PM2.5 vs CHOLSCREEN_CrudePrev for 2022



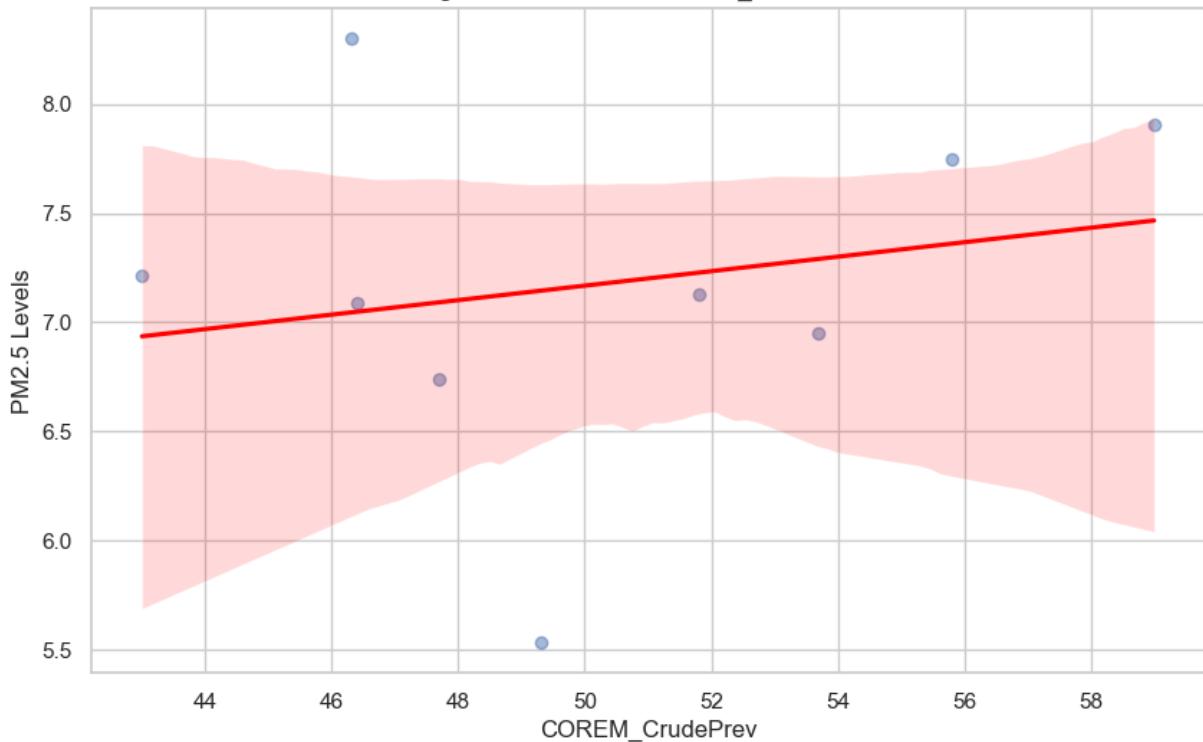
Linear Regression: PM2.5 vs COLON_SCREEN_CrudePrev for 2022



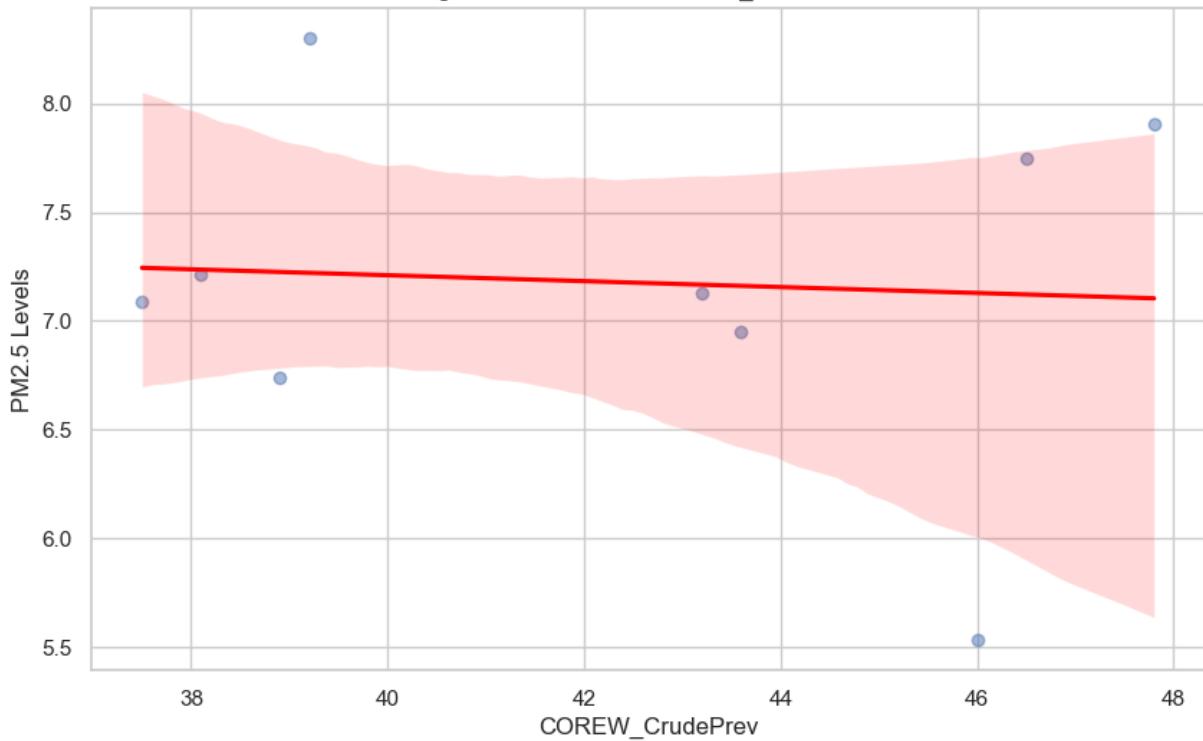
Linear Regression: PM2.5 vs COPD_CrudePrev for 2022

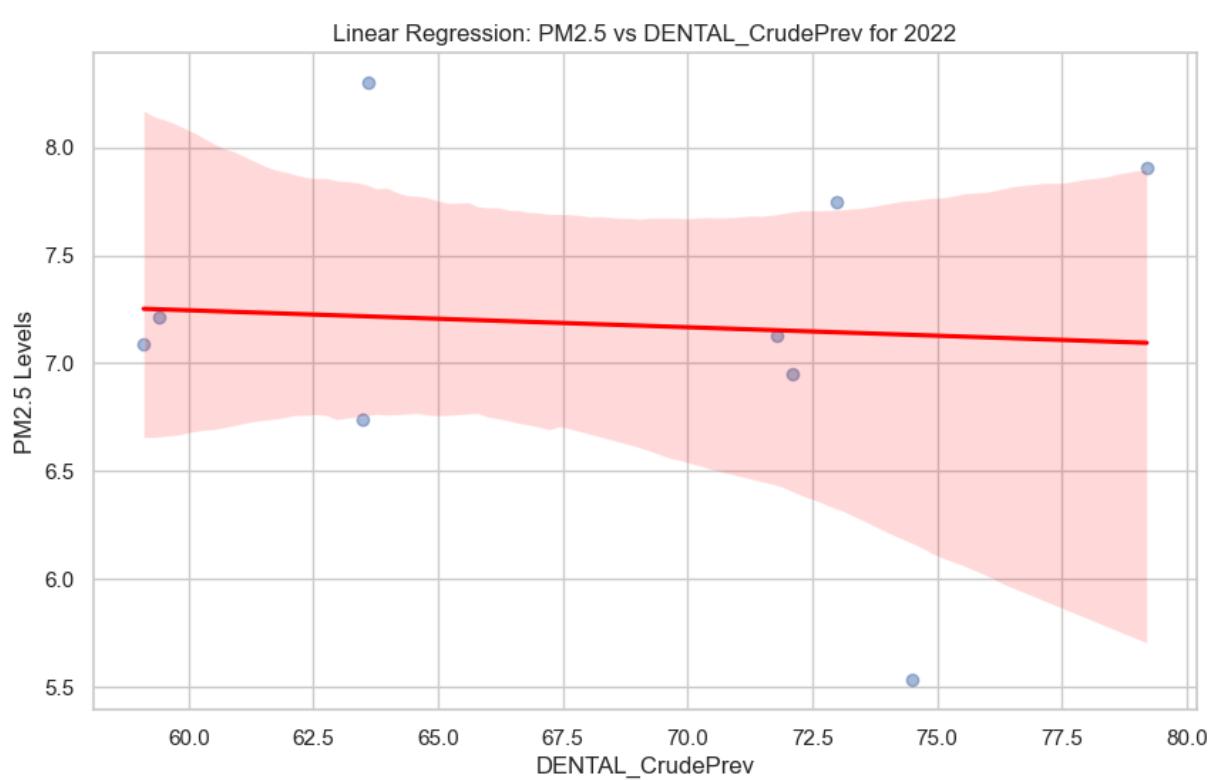
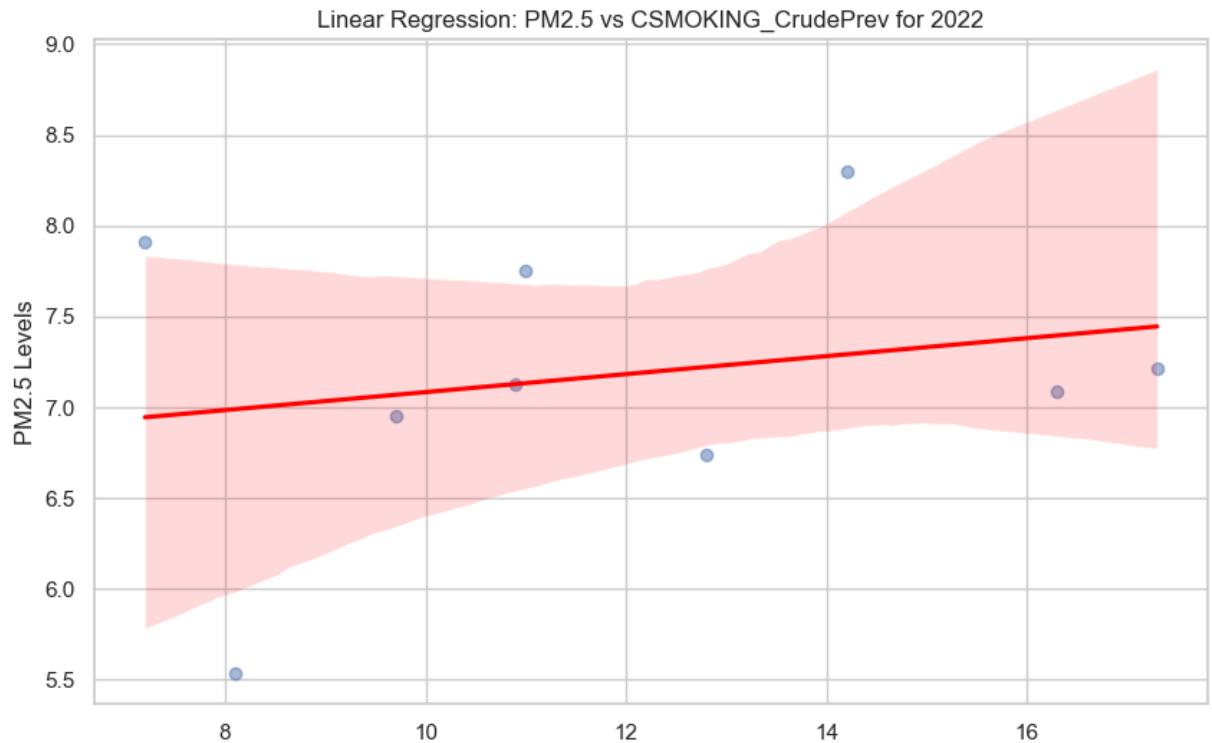


Linear Regression: PM2.5 vs COREM_CrudePrev for 2022

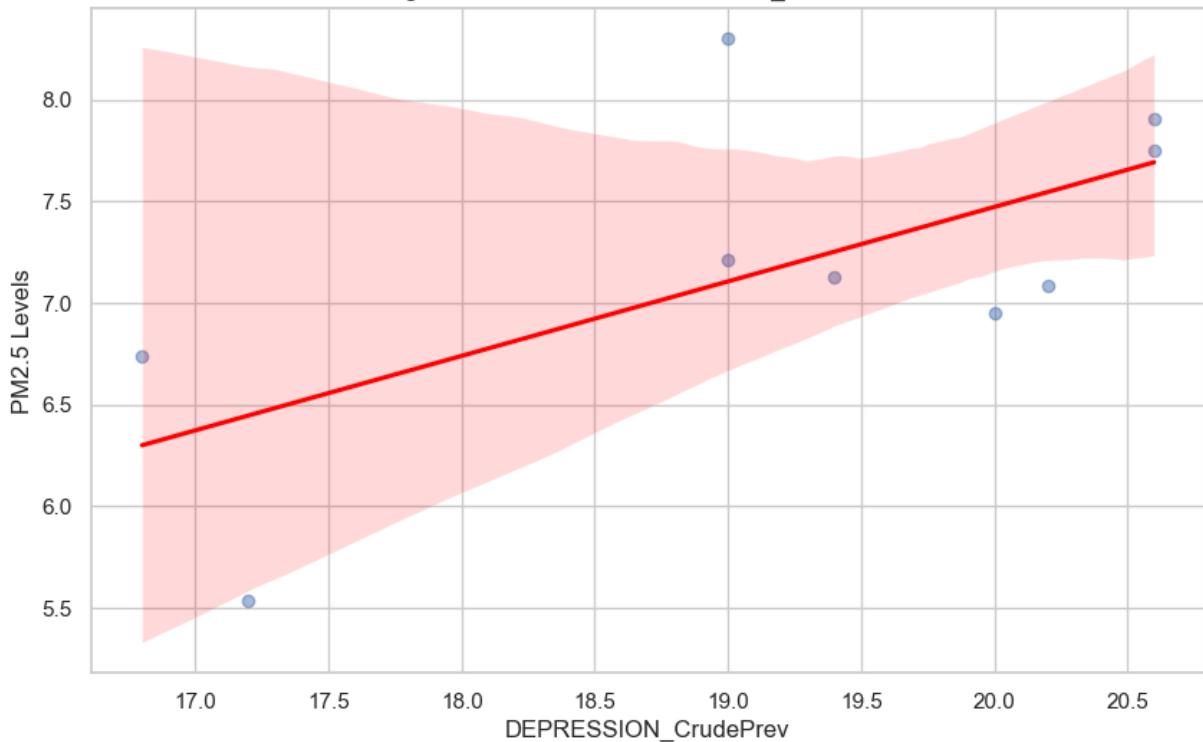


Linear Regression: PM2.5 vs COREW_CrudePrev for 2022

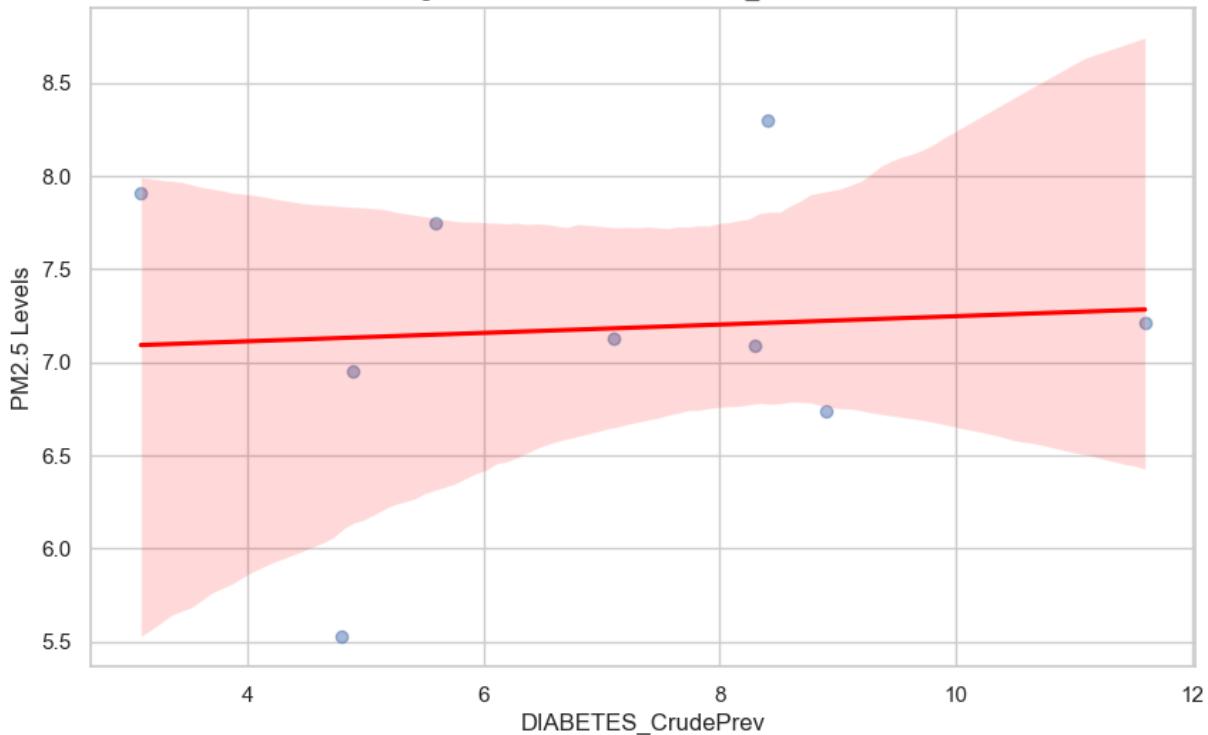


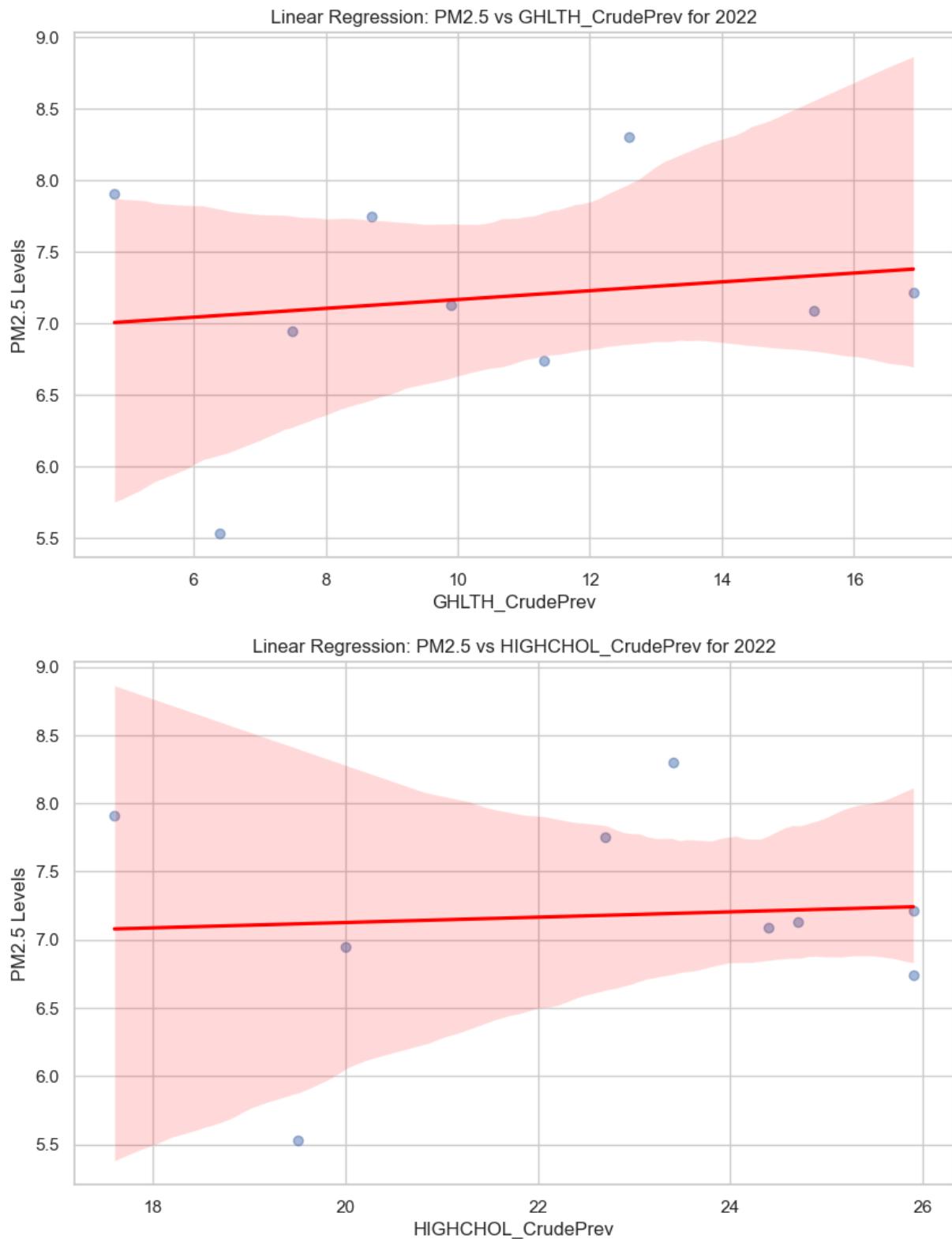


Linear Regression: PM2.5 vs DEPRESSION_CrudePrev for 2022

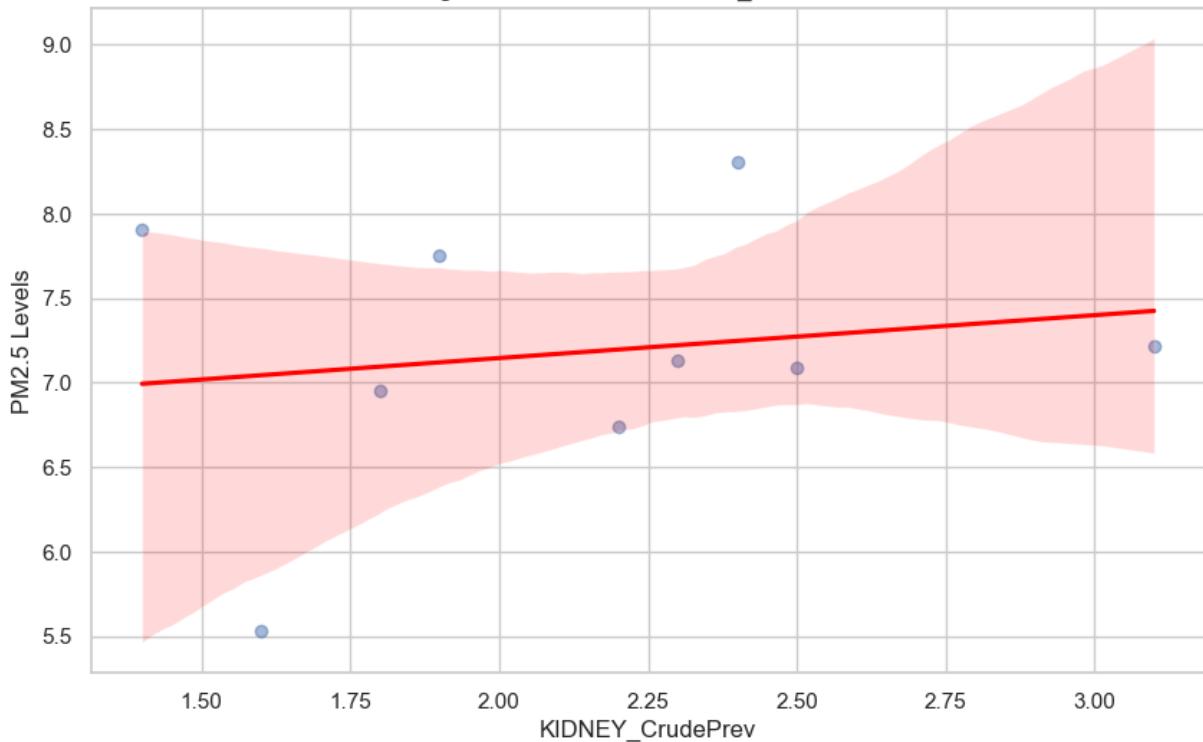


Linear Regression: PM2.5 vs DIABETES_CrudePrev for 2022

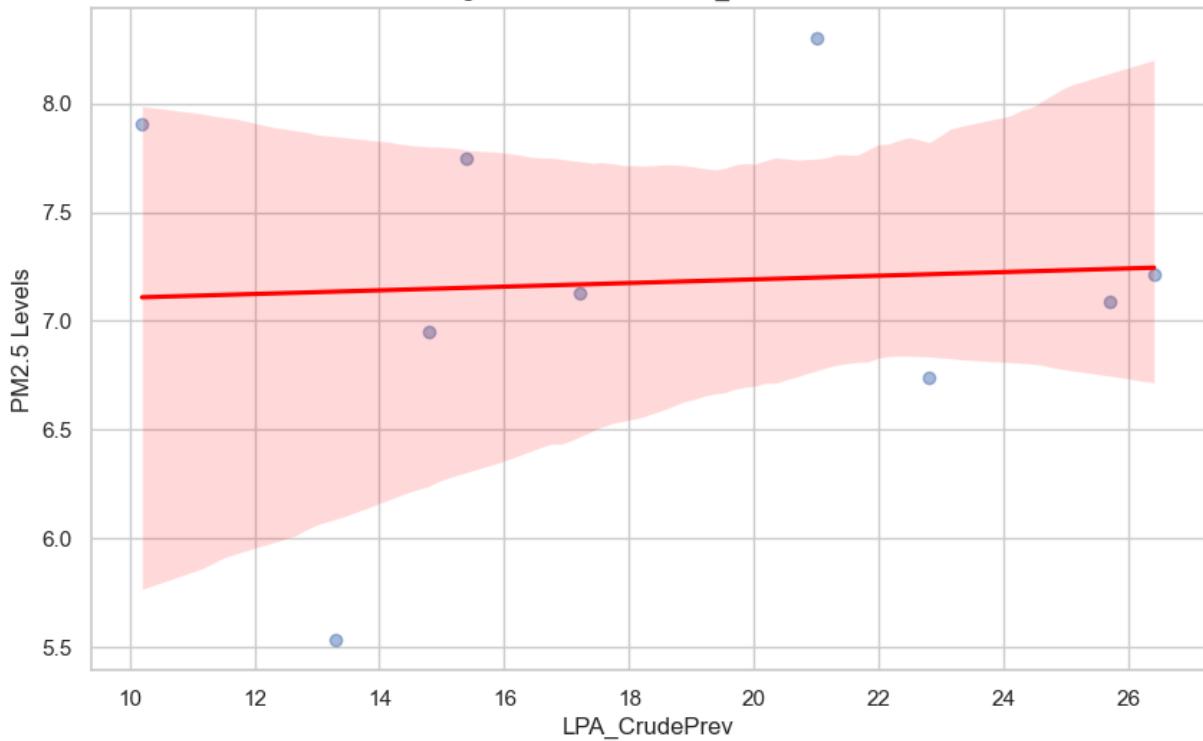




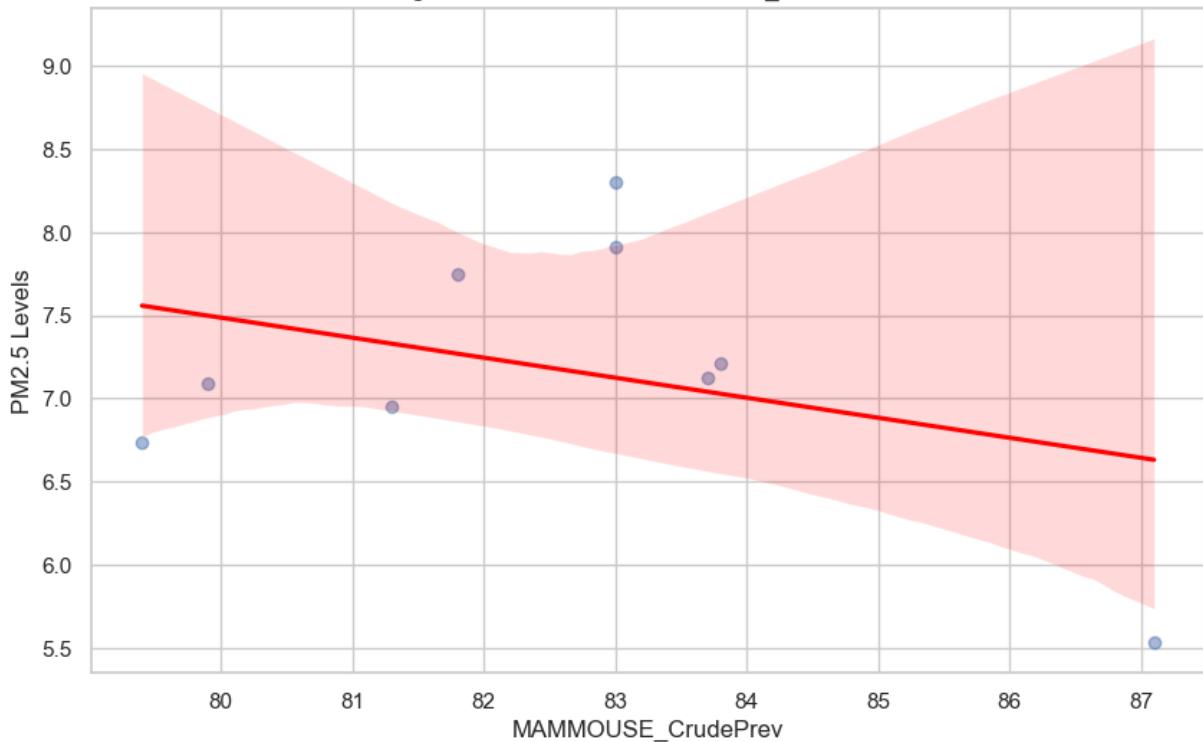
Linear Regression: PM2.5 vs KIDNEY_CrudePrev for 2022



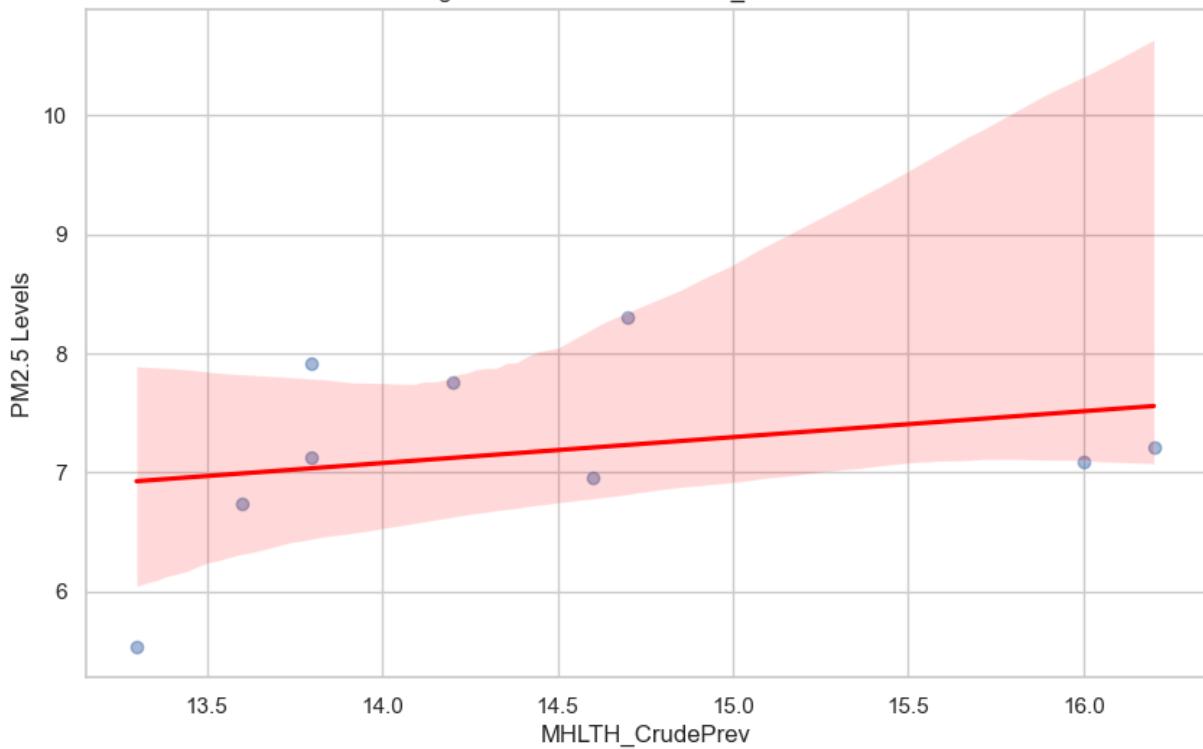
Linear Regression: PM2.5 vs LPA_CrudePrev for 2022



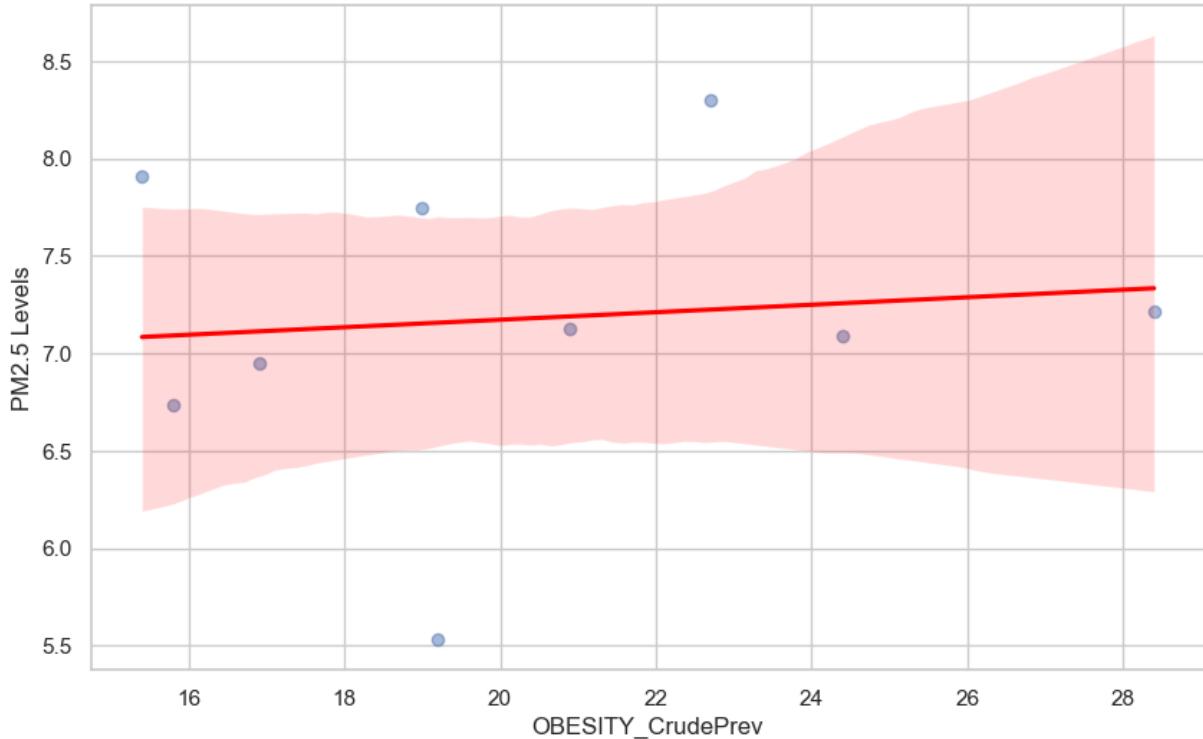
Linear Regression: PM2.5 vs MAMMOUSE_CrudePrev for 2022



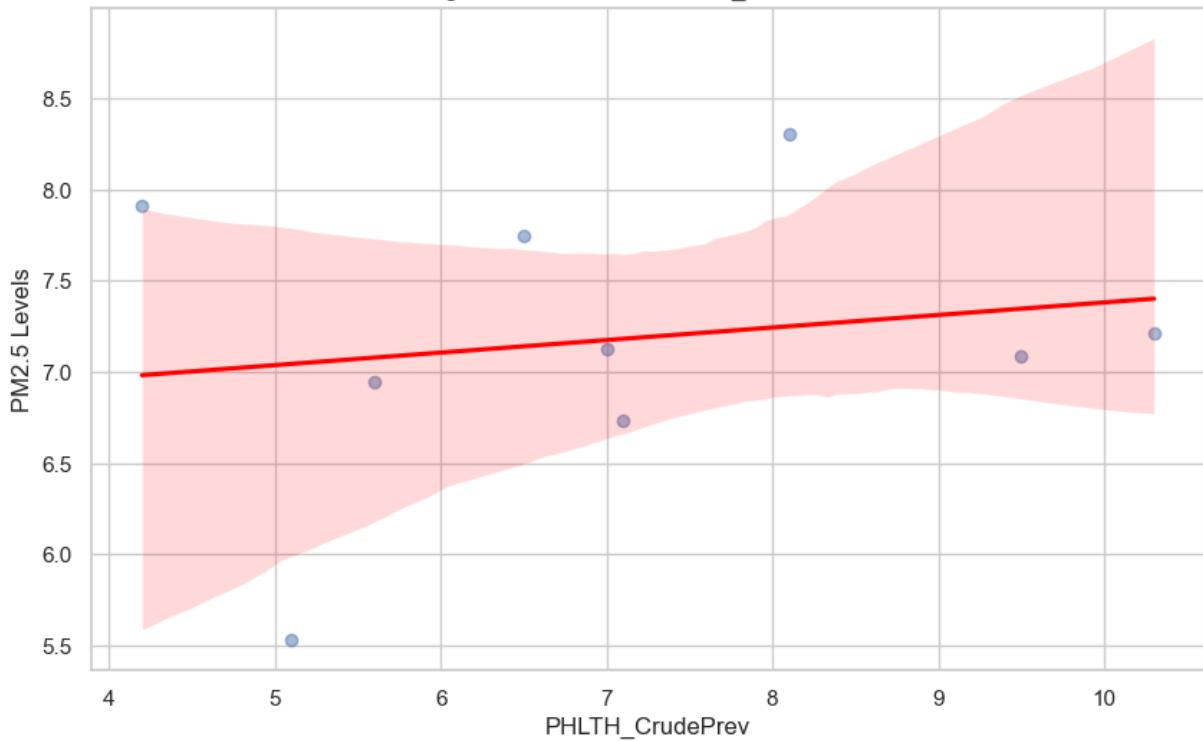
Linear Regression: PM2.5 vs MHLTH_CrudePrev for 2022



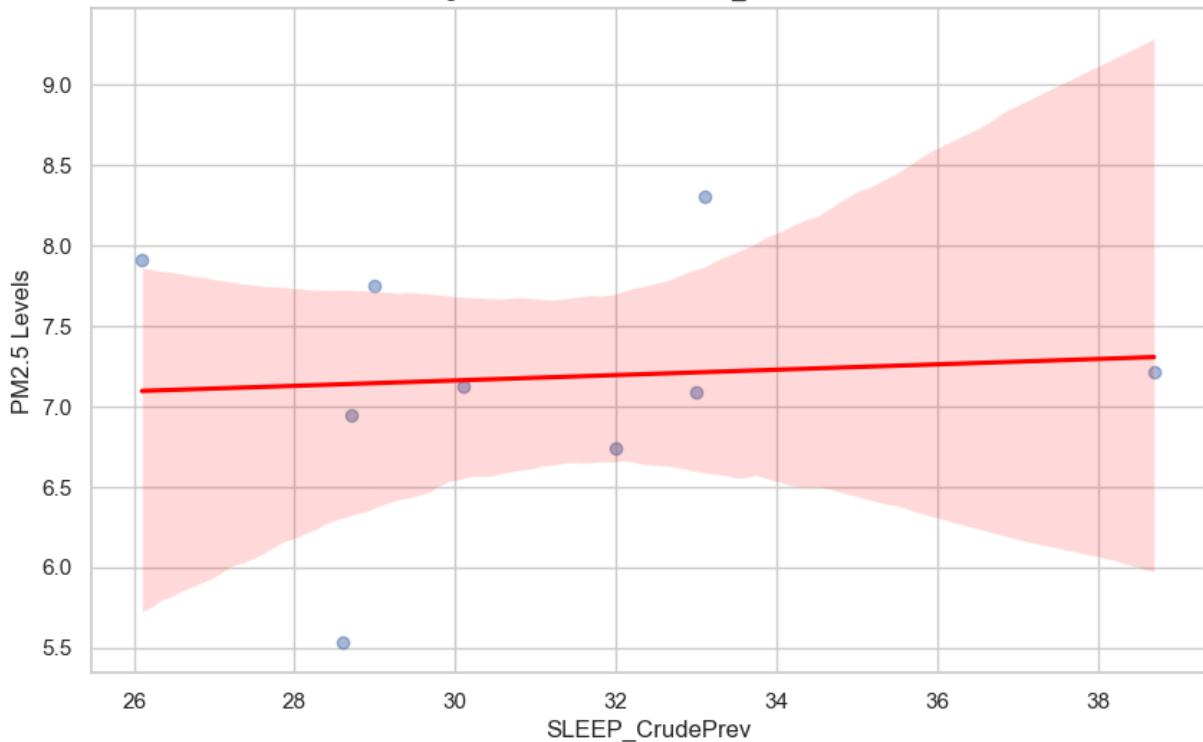
Linear Regression: PM2.5 vs OBESITY_CrudePrev for 2022



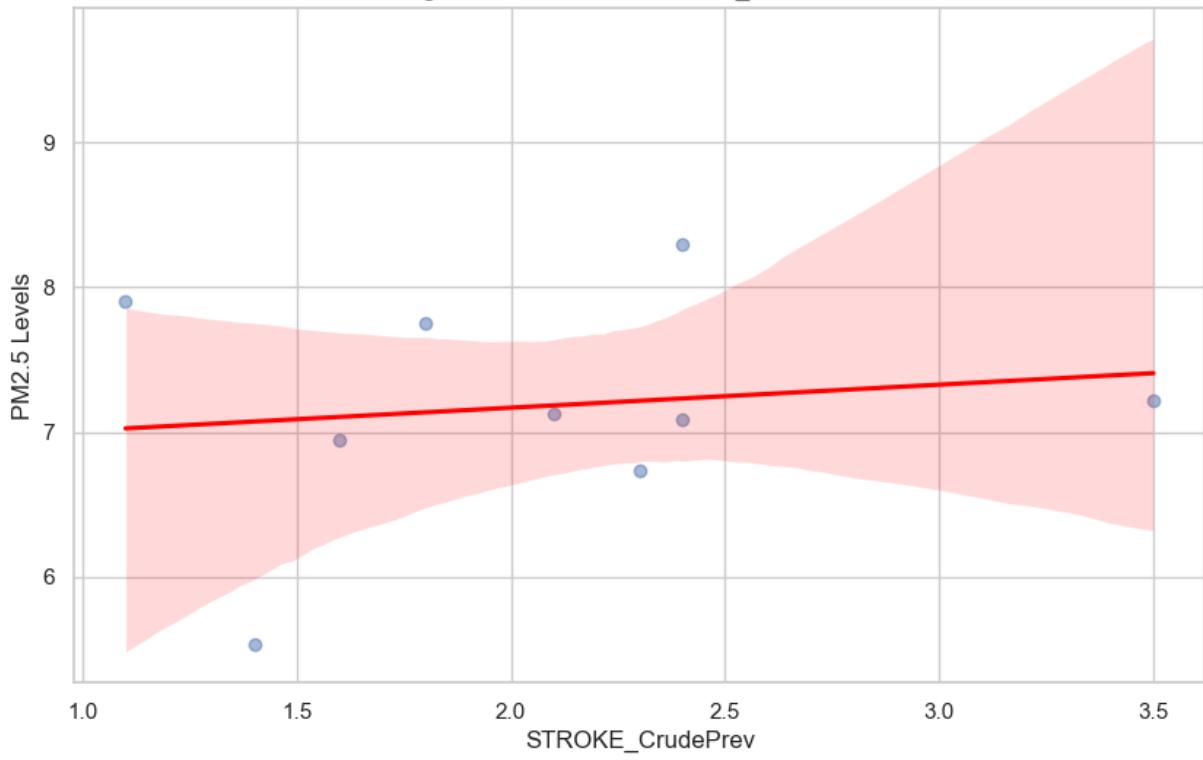
Linear Regression: PM2.5 vs PHLTH_CrudePrev for 2022

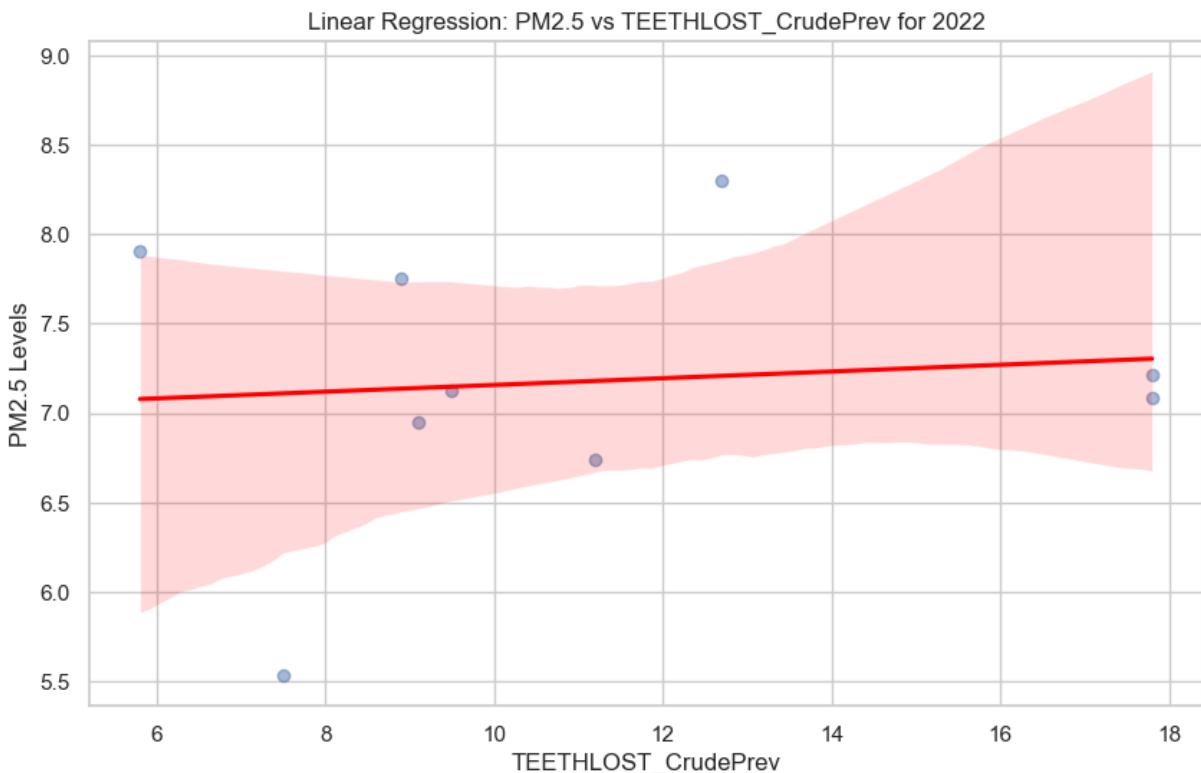


Linear Regression: PM2.5 vs SLEEP_CrudePrev for 2022



Linear Regression: PM2.5 vs STROKE_CrudePrev for 2022

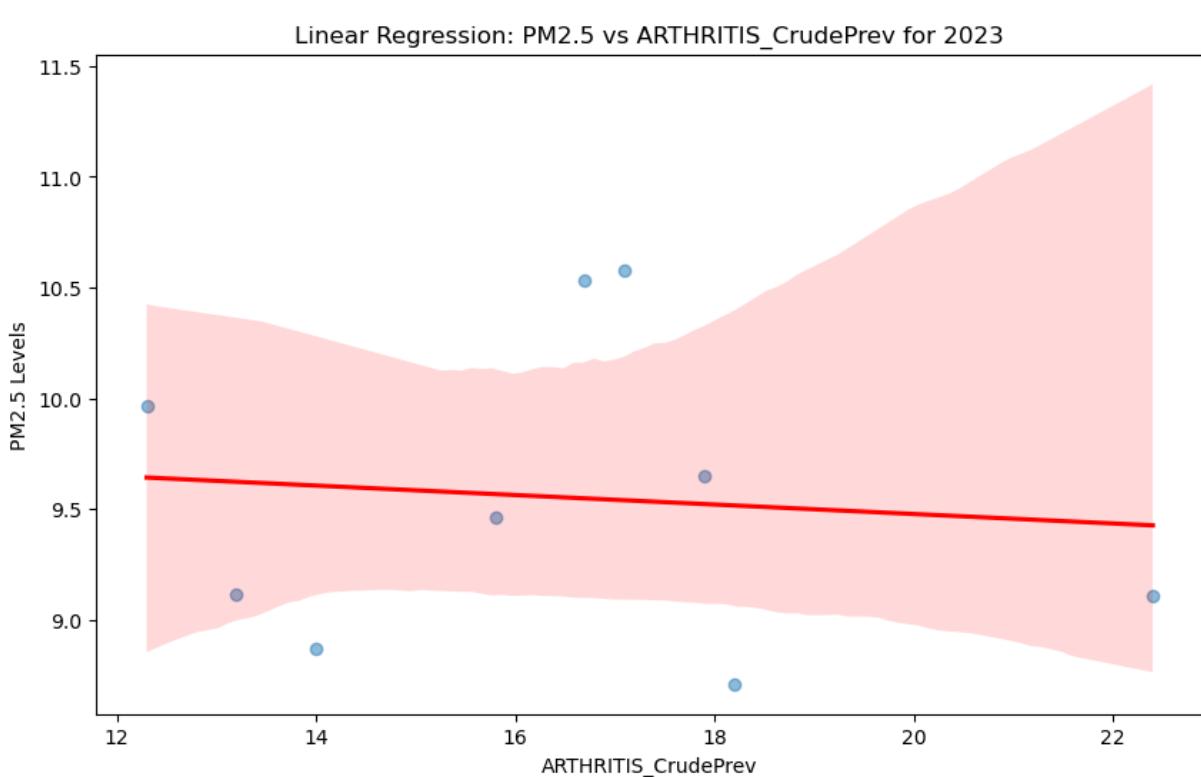
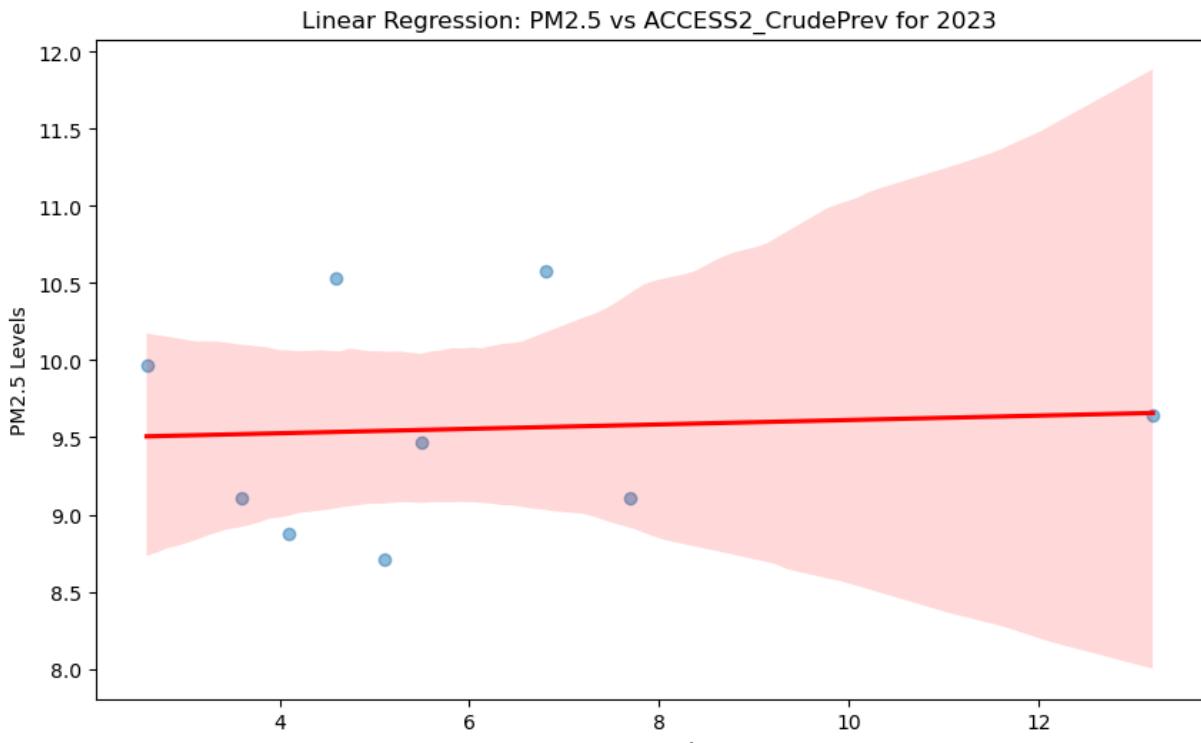


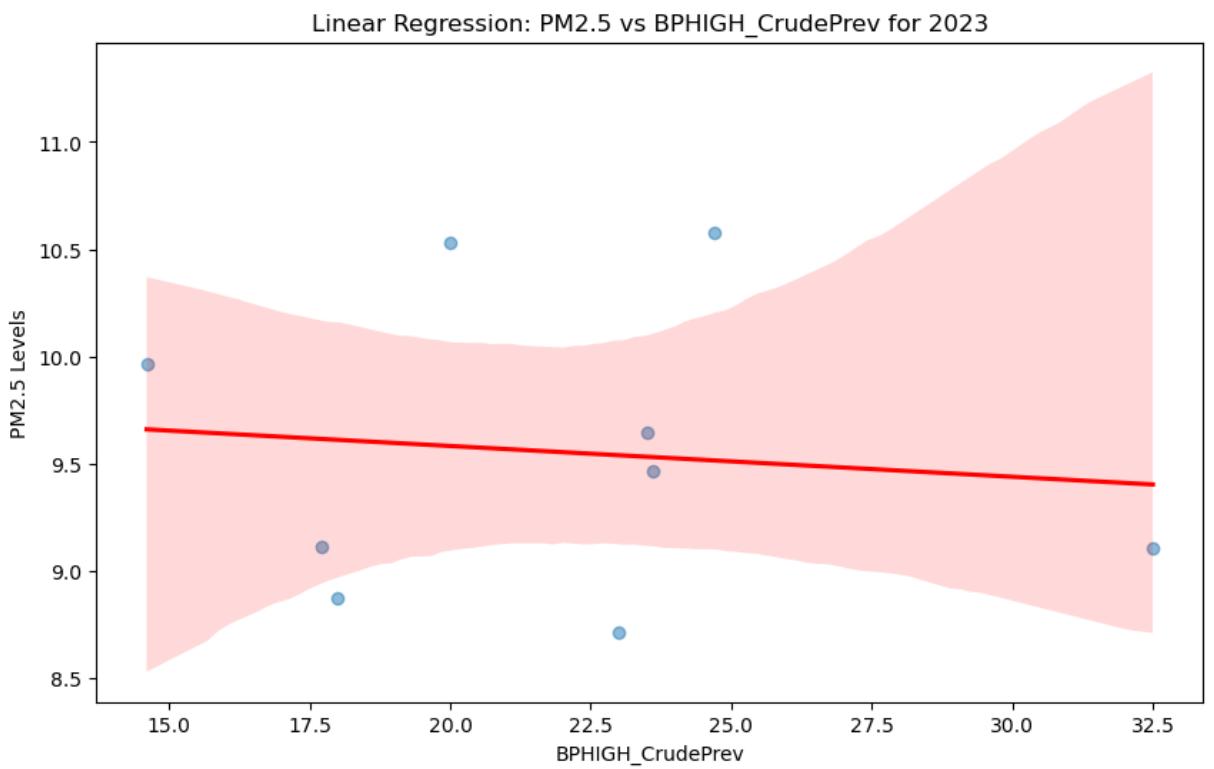
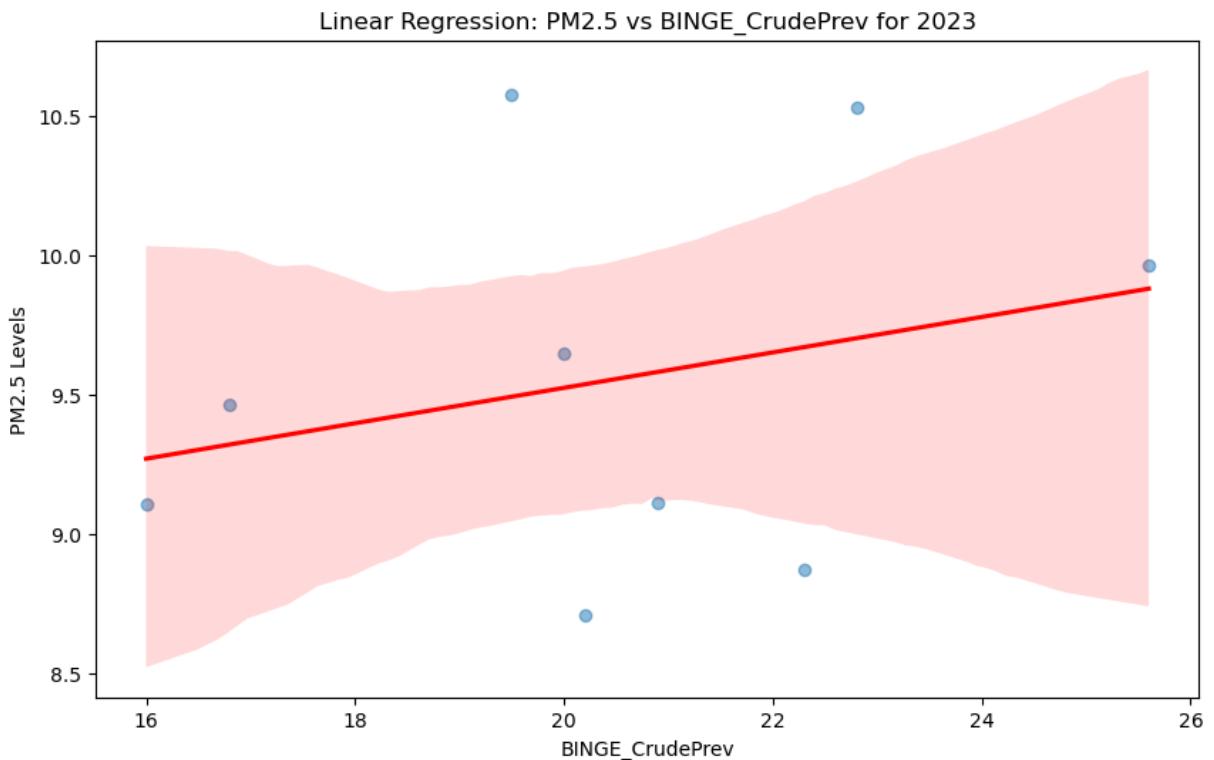


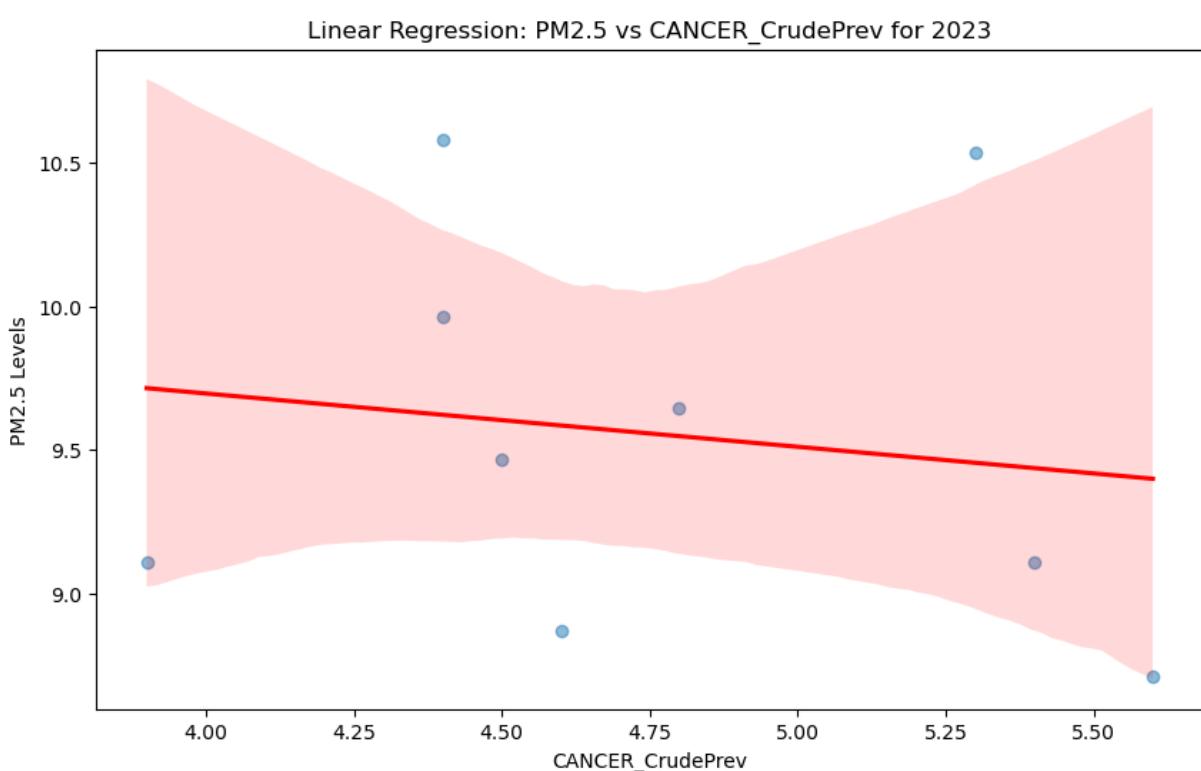
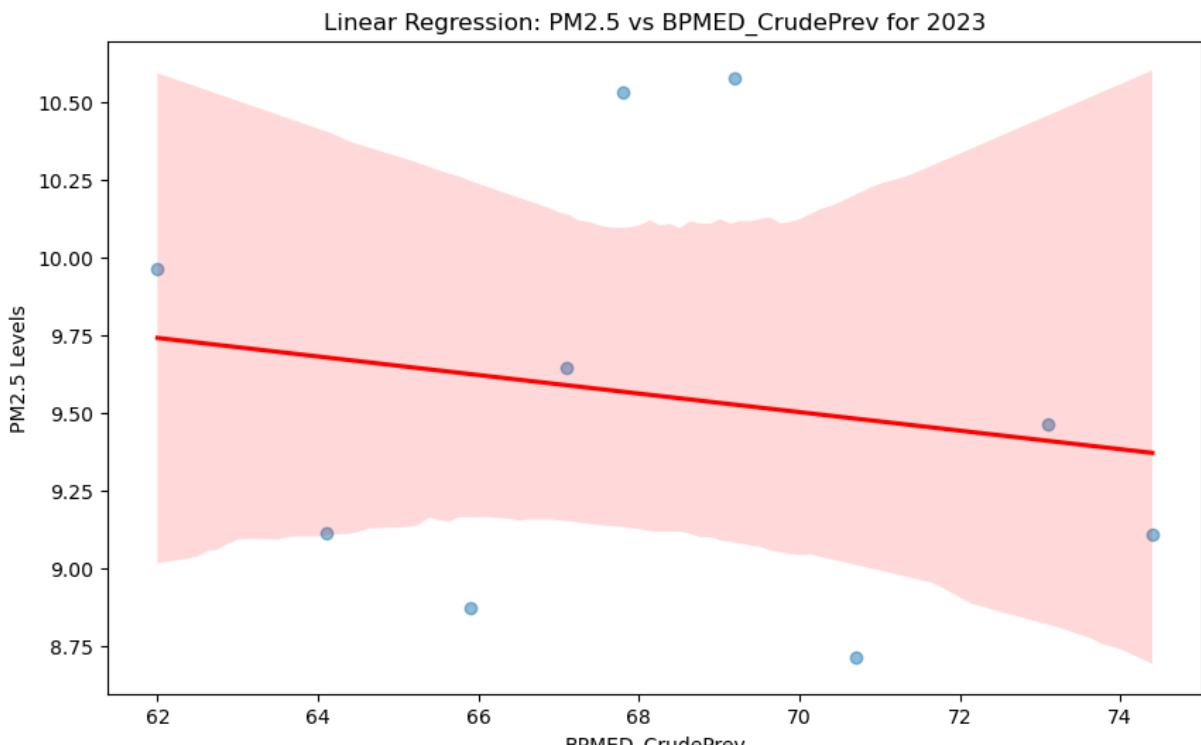
```
In [10]: def plot_linear_regression(data, crude_columns, year):
    for column in crude_columns:
        if column != 'PM2.5':
            # Performing linear regression
            slope, intercept, r_value, p_value, std_err = stats.linregress(column, data['PM2.5'])

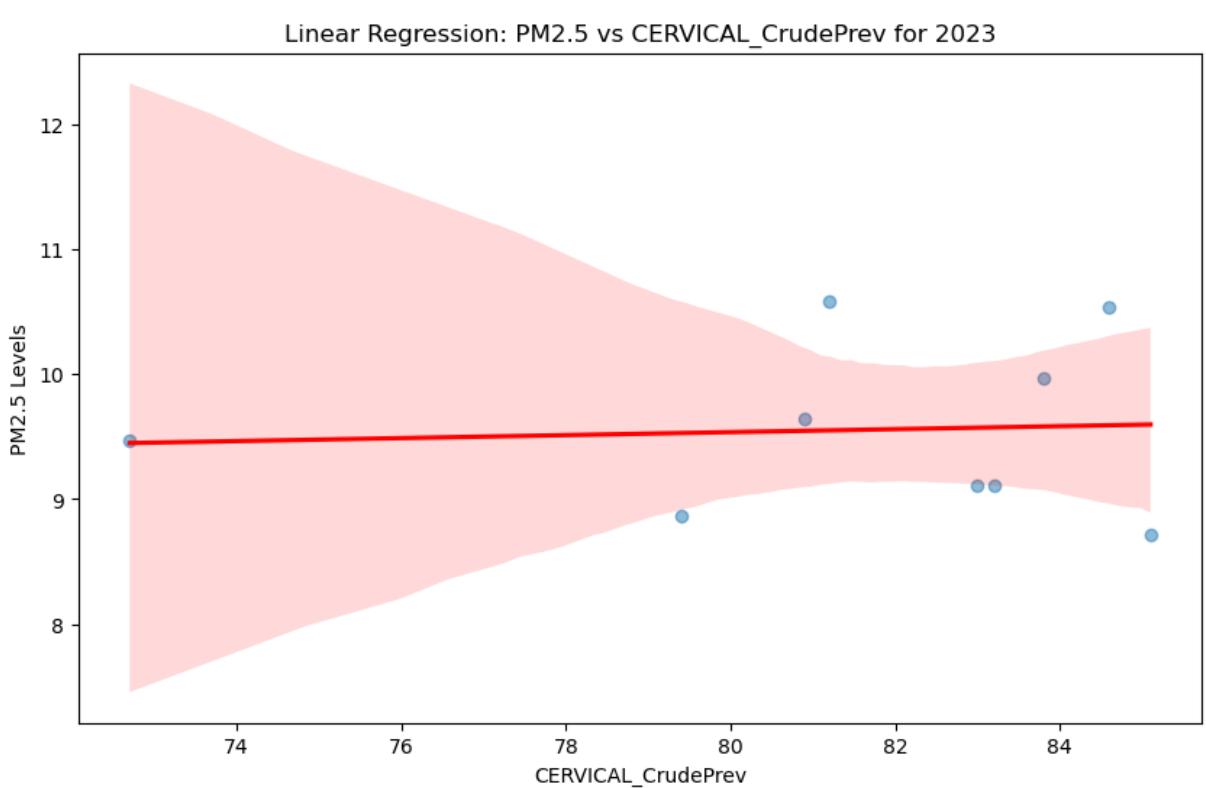
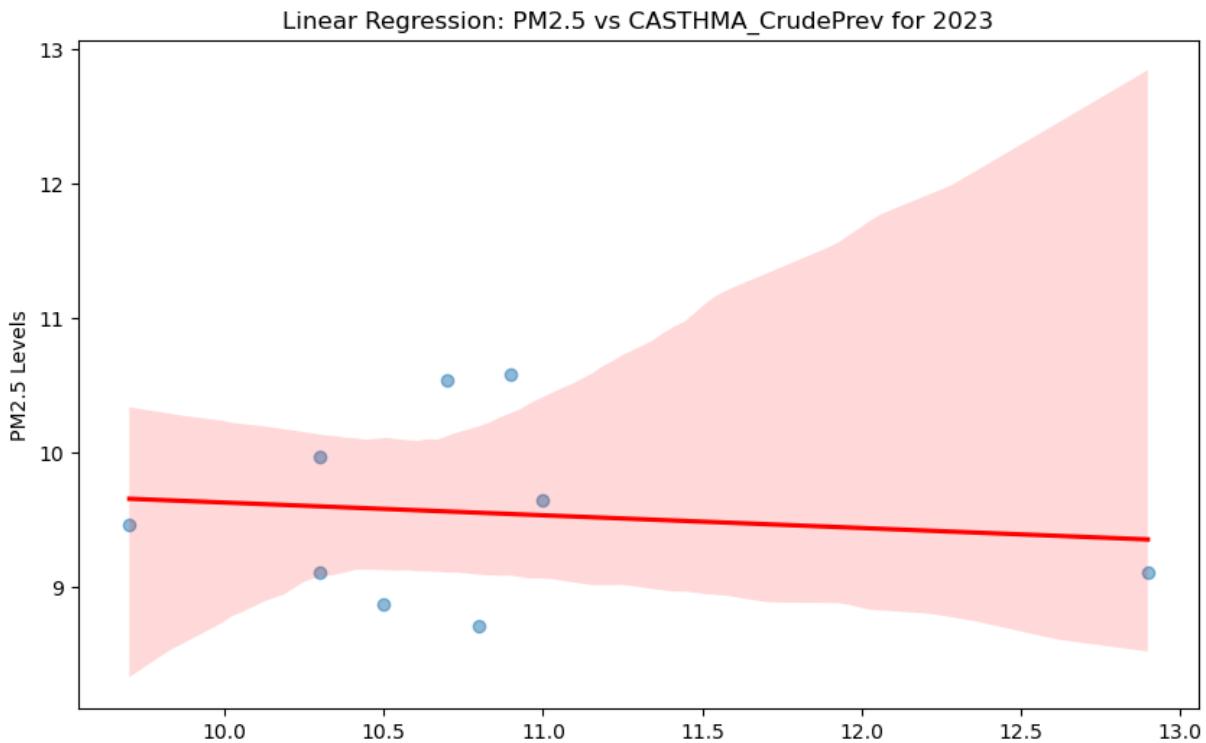
            # Plotting
            plt.figure(figsize=(10, 6))
            sns.regplot(x=column, y='PM2.5', data=data, scatter_kws={'alpha': 0.5})
            plt.title(f'Linear Regression: PM2.5 vs {column} for {year}')
            plt.xlabel(column)
            plt.ylabel('PM2.5 Levels')
            plt.show()

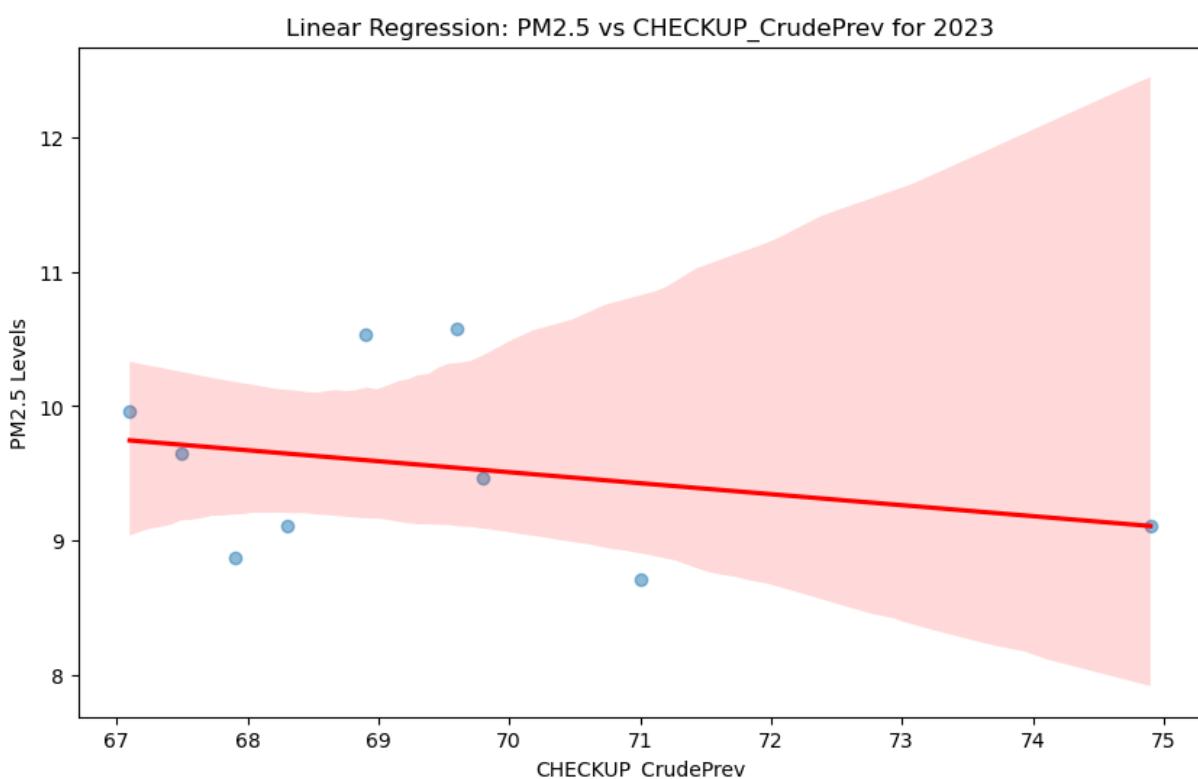
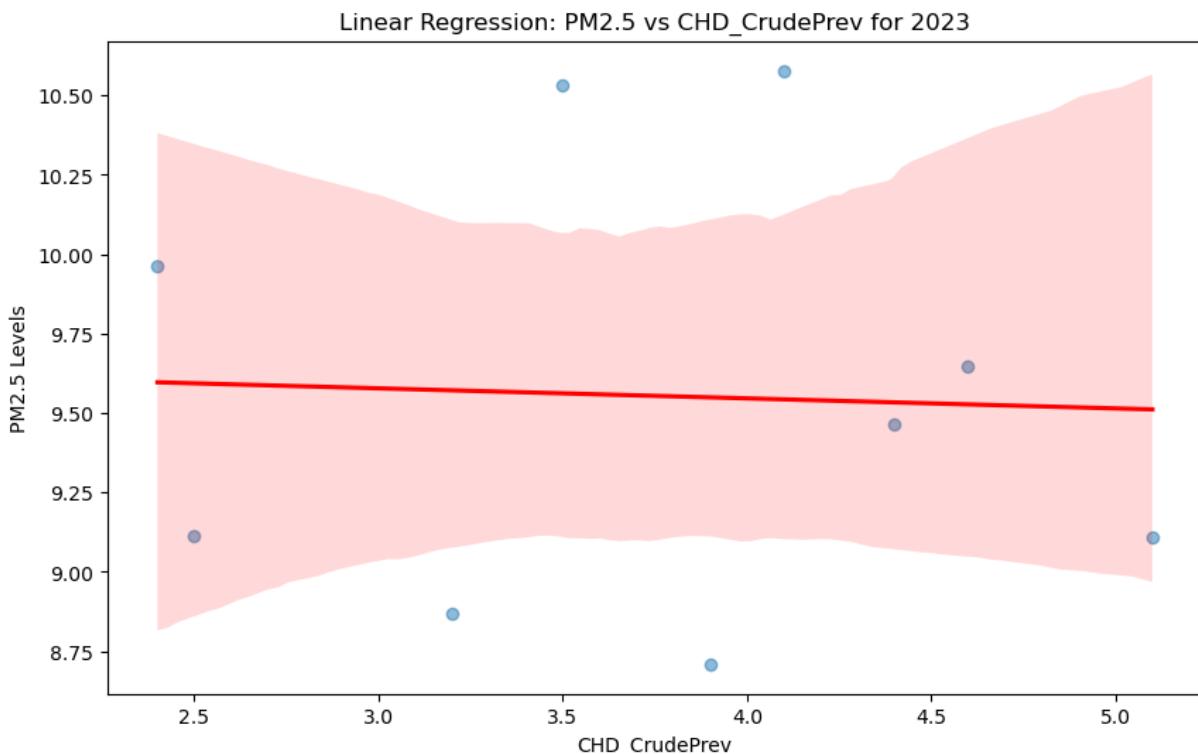
crude_columns_2023 = [col for col in health_data_2023.columns if 'Crude' in col]
plot_linear_regression(health_data_2023, crude_columns_2023, '2023')
```

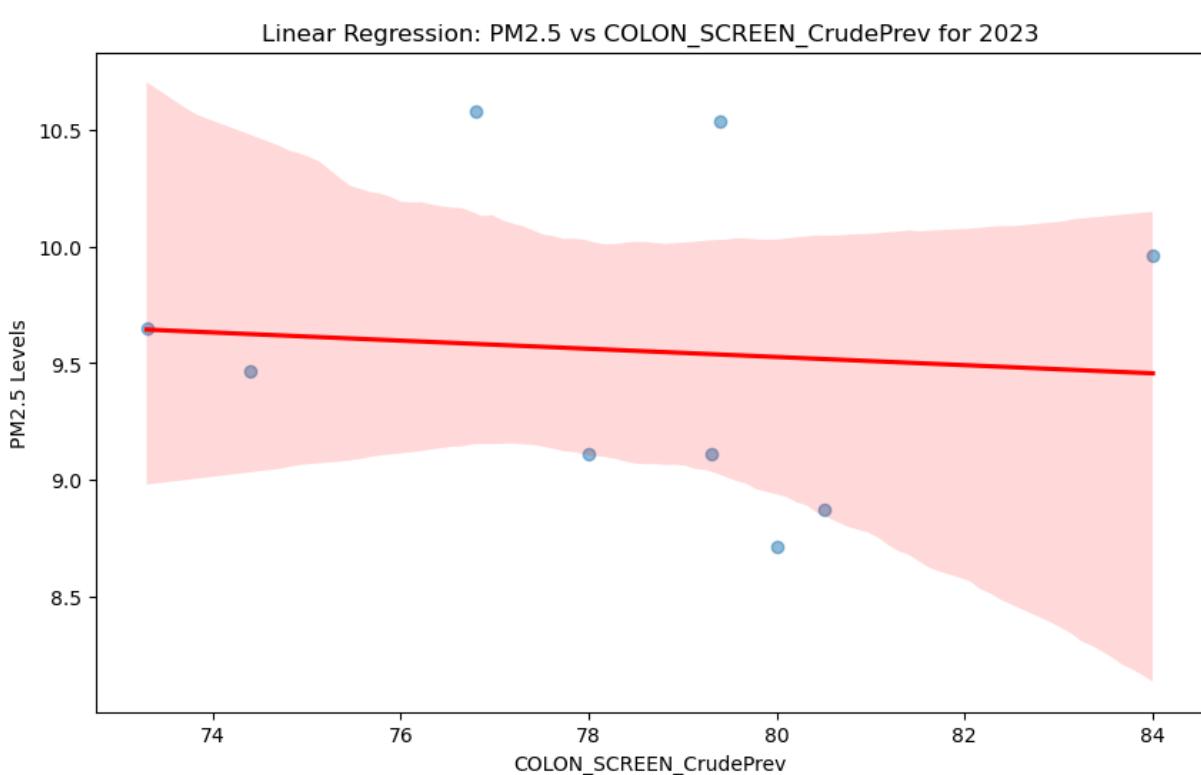
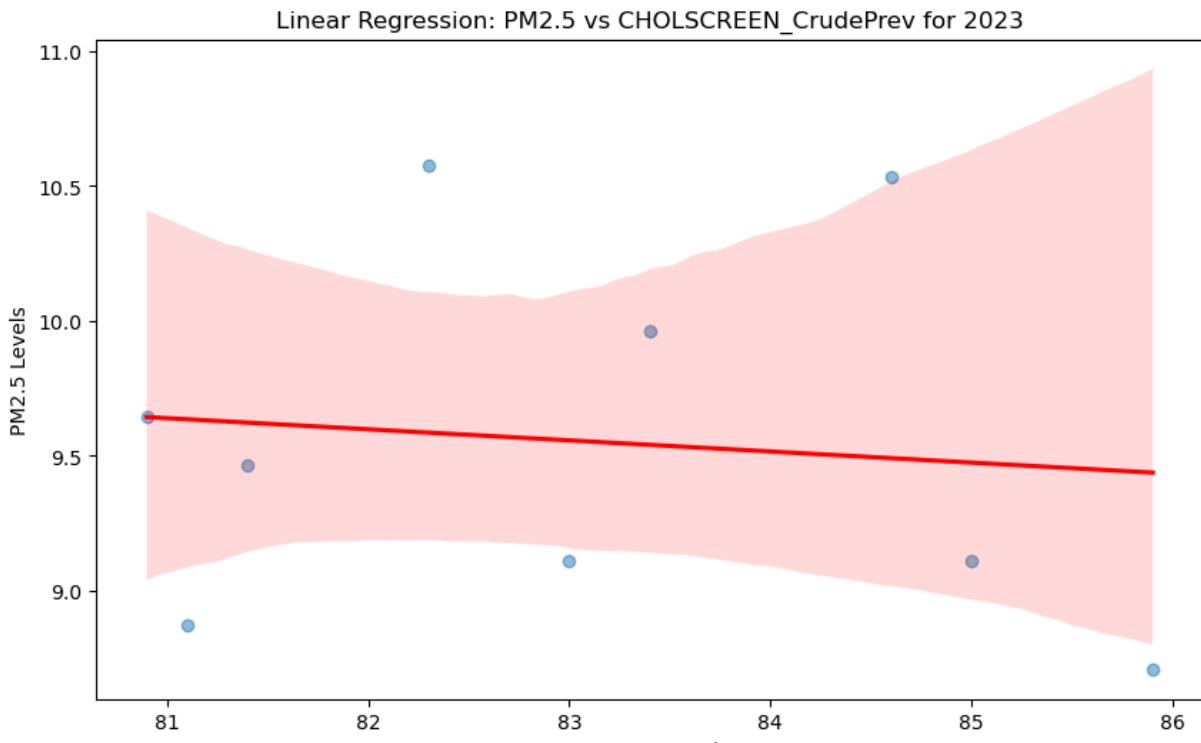


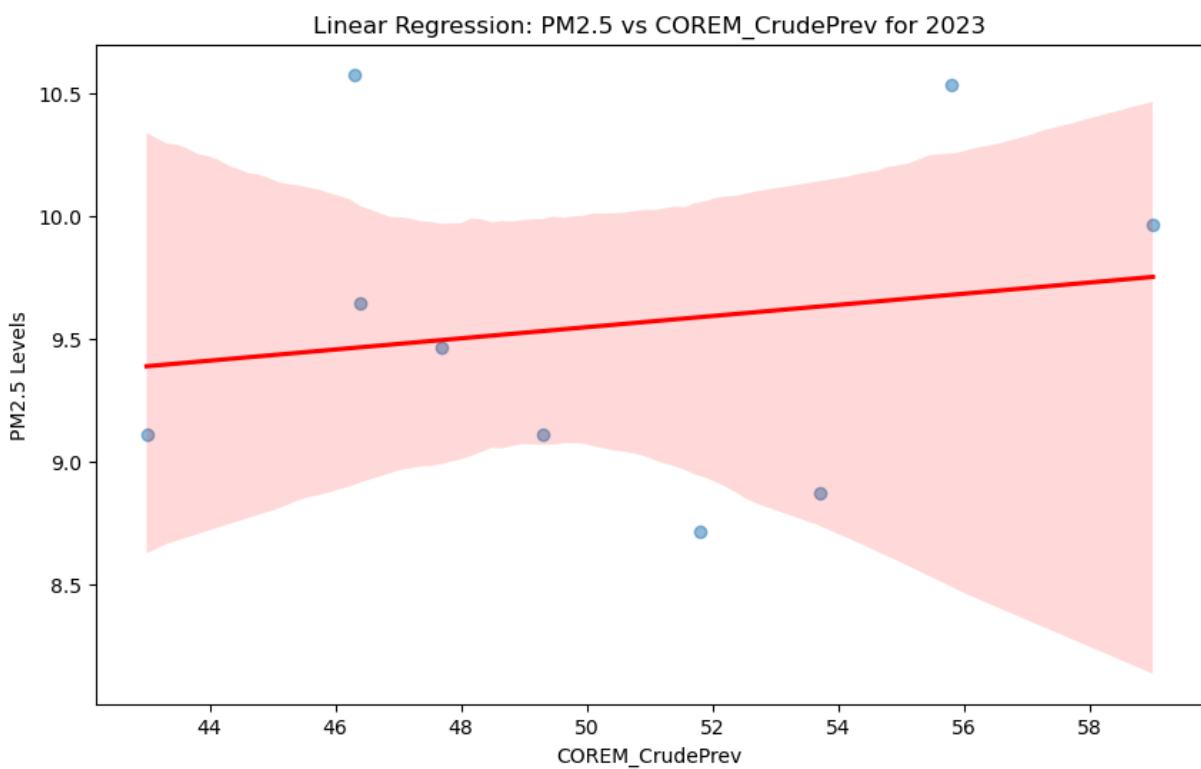
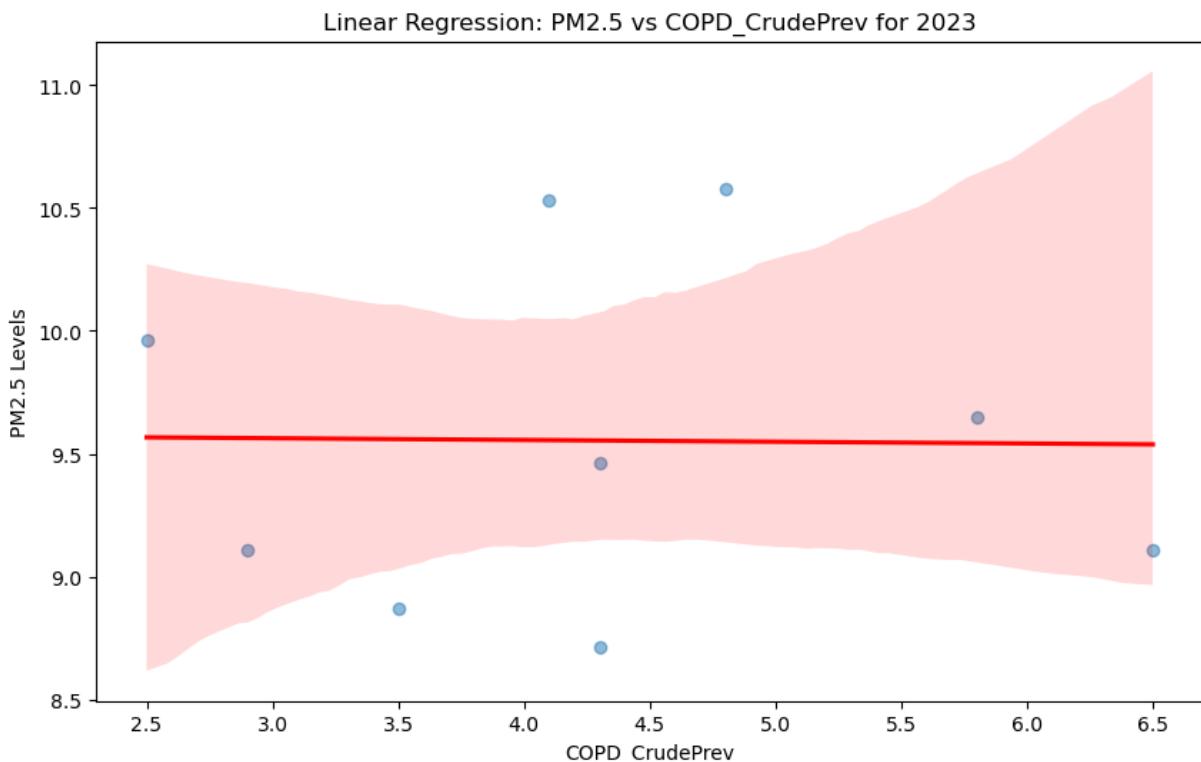


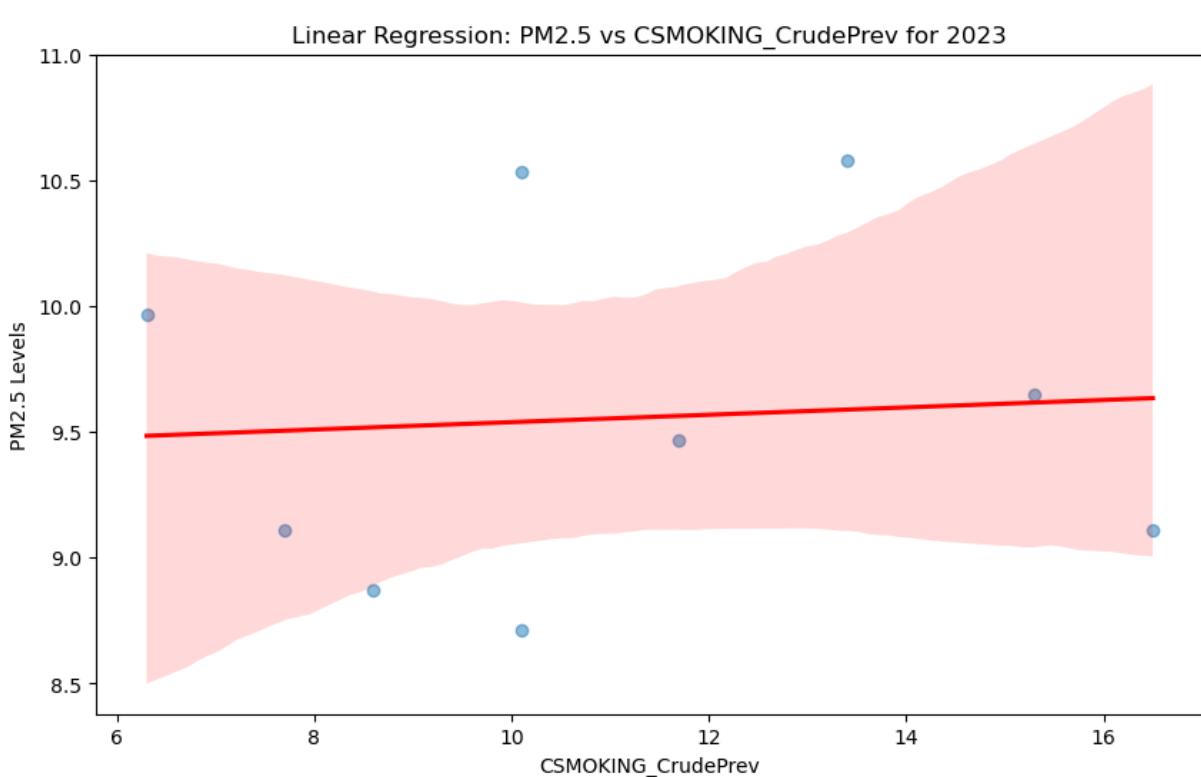
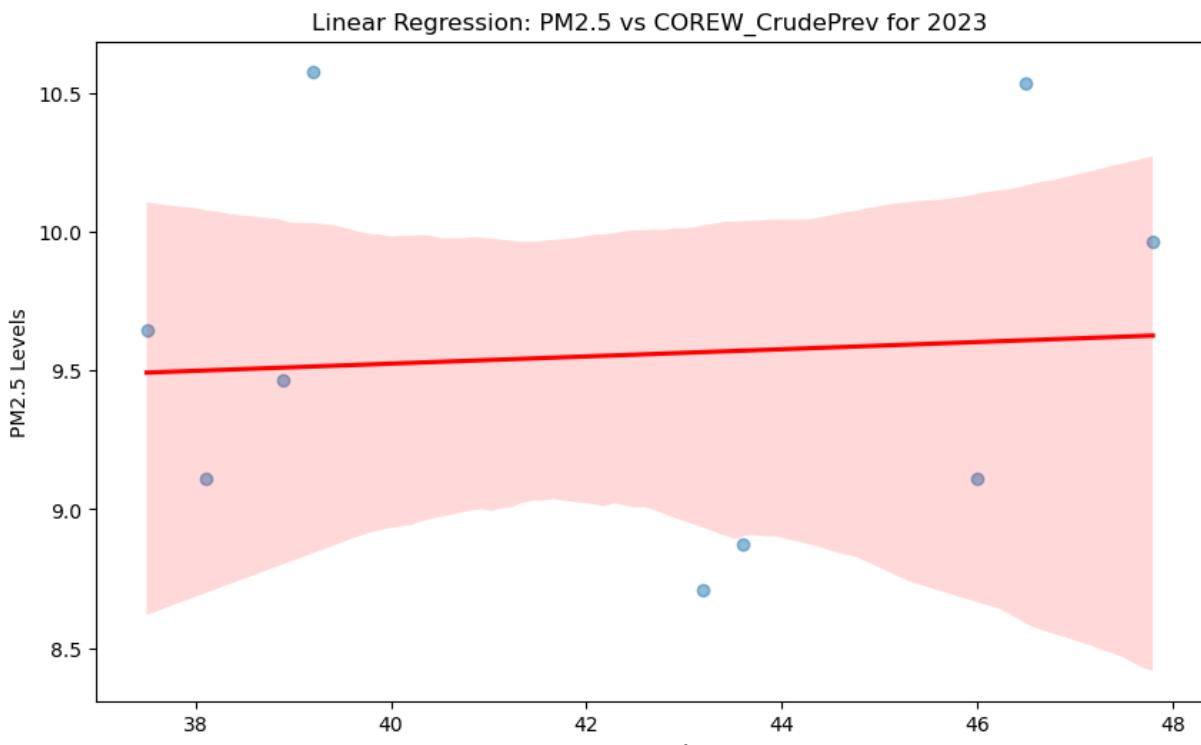


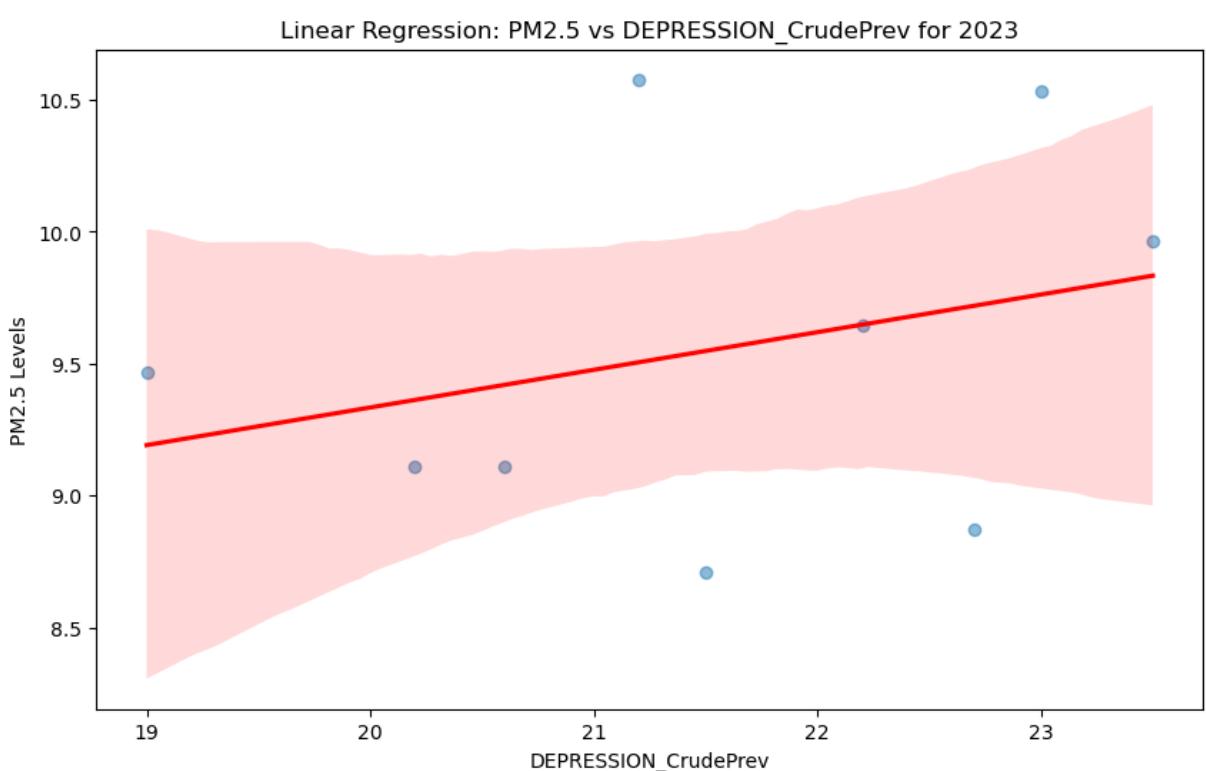
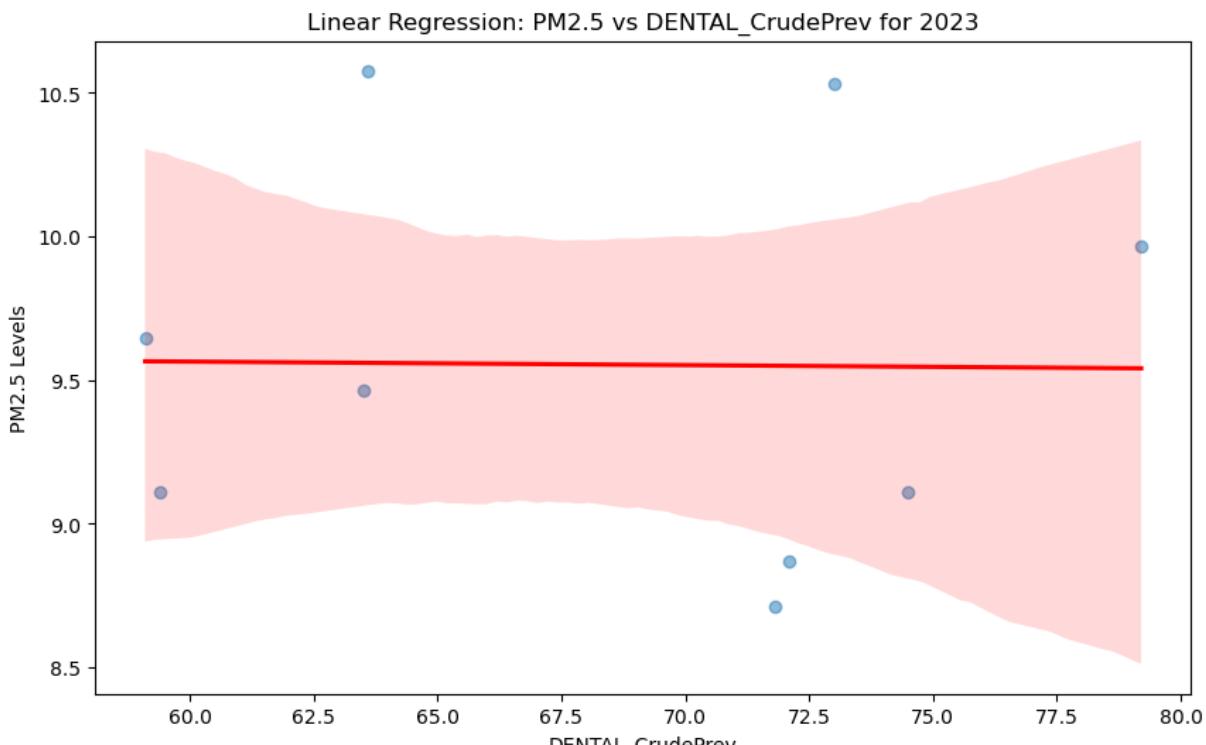




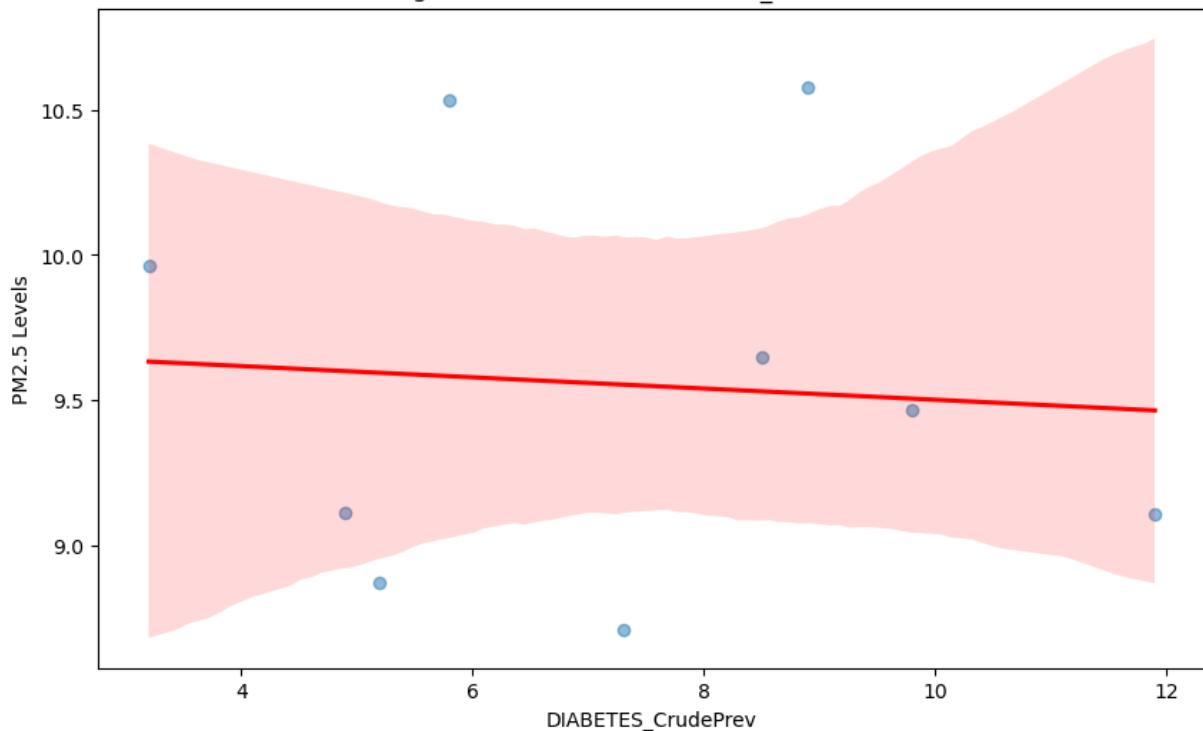




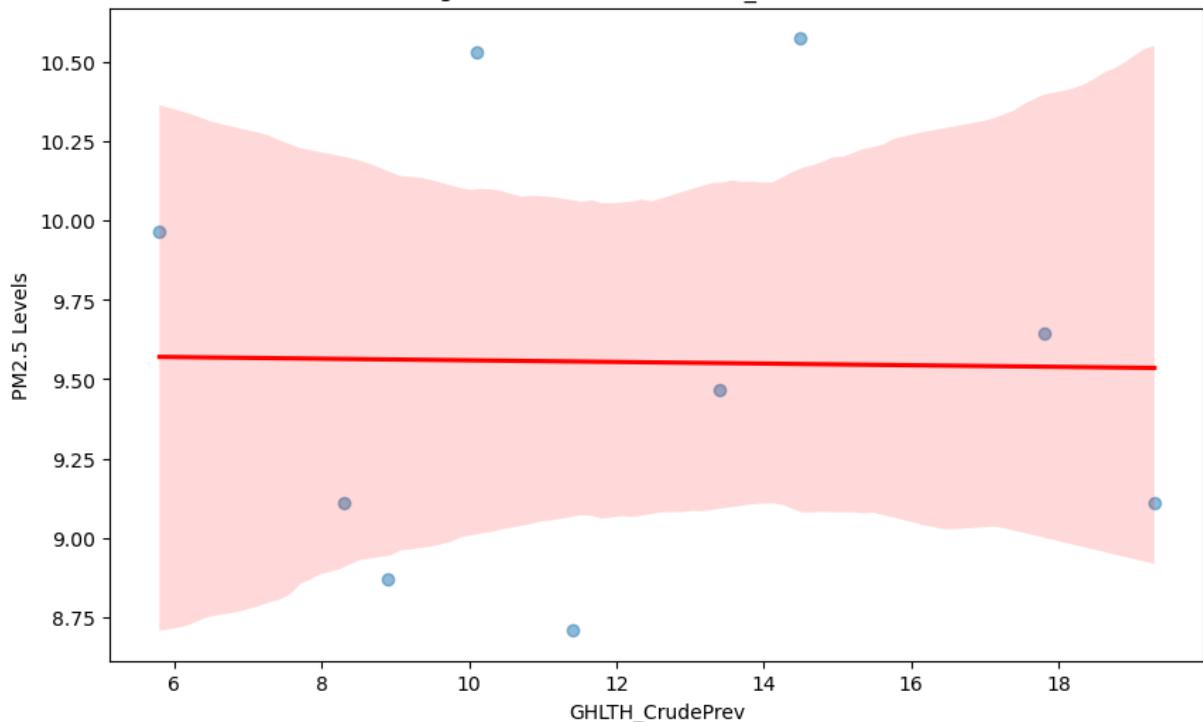


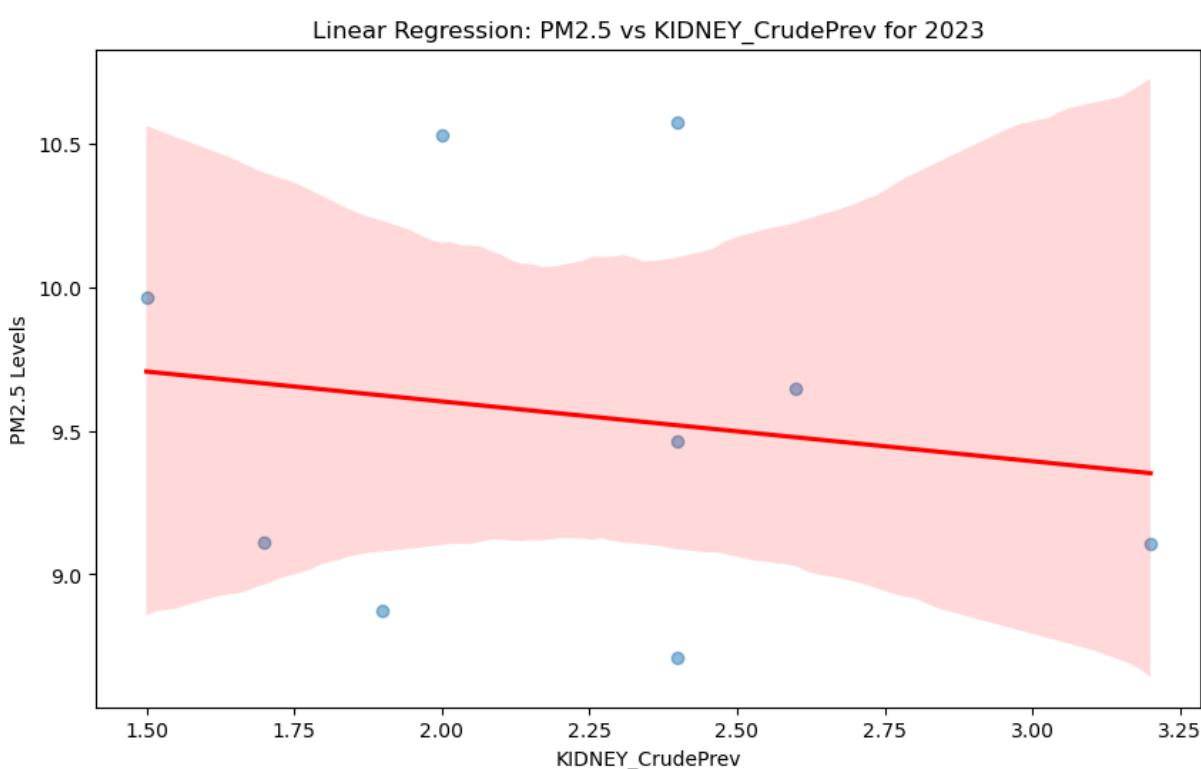
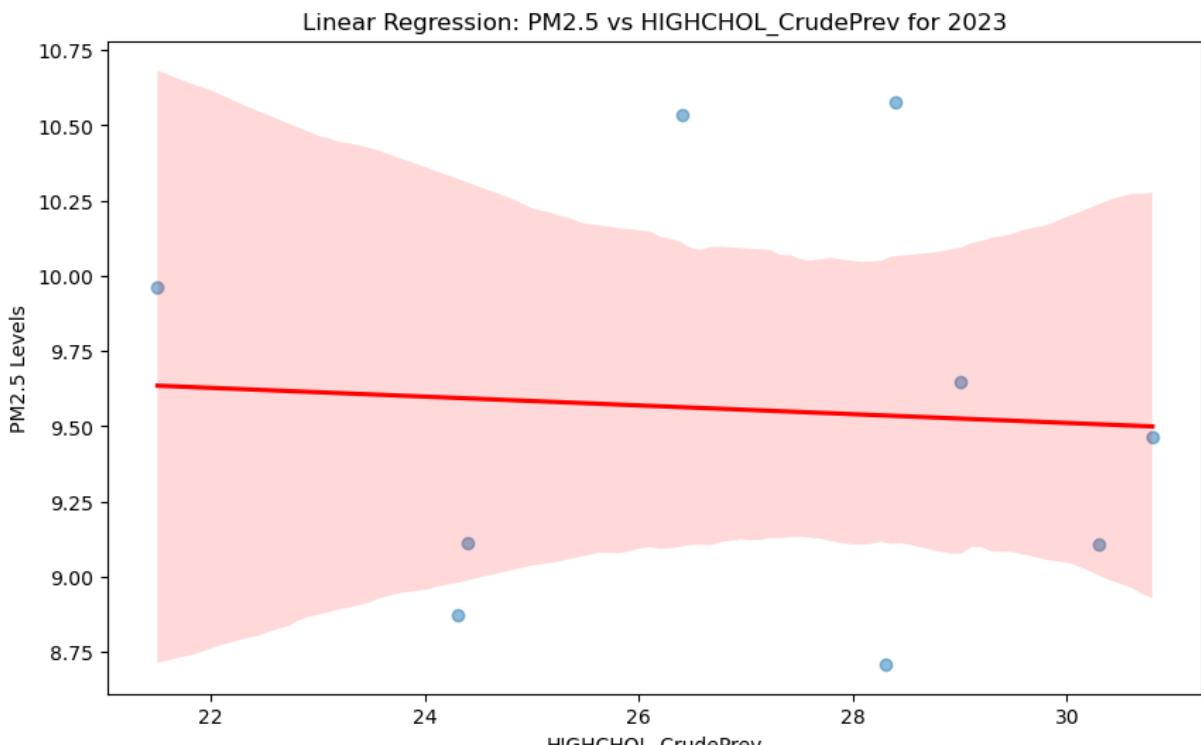


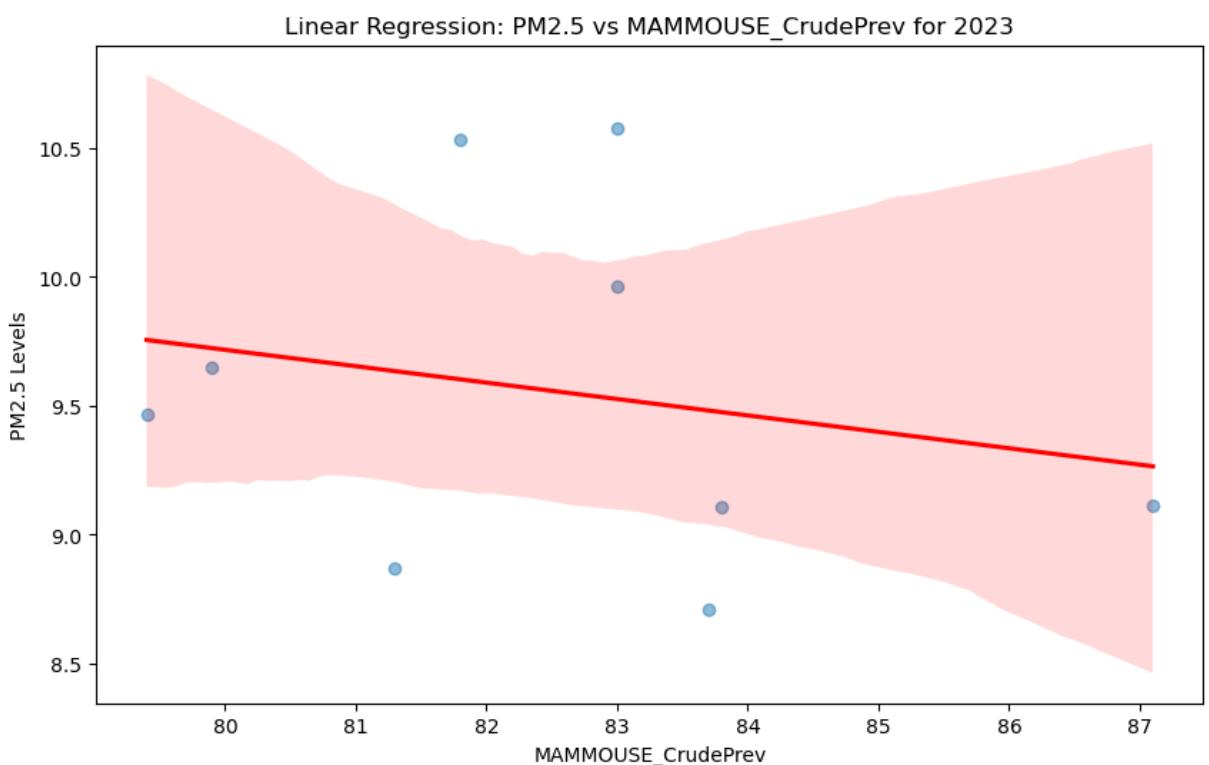
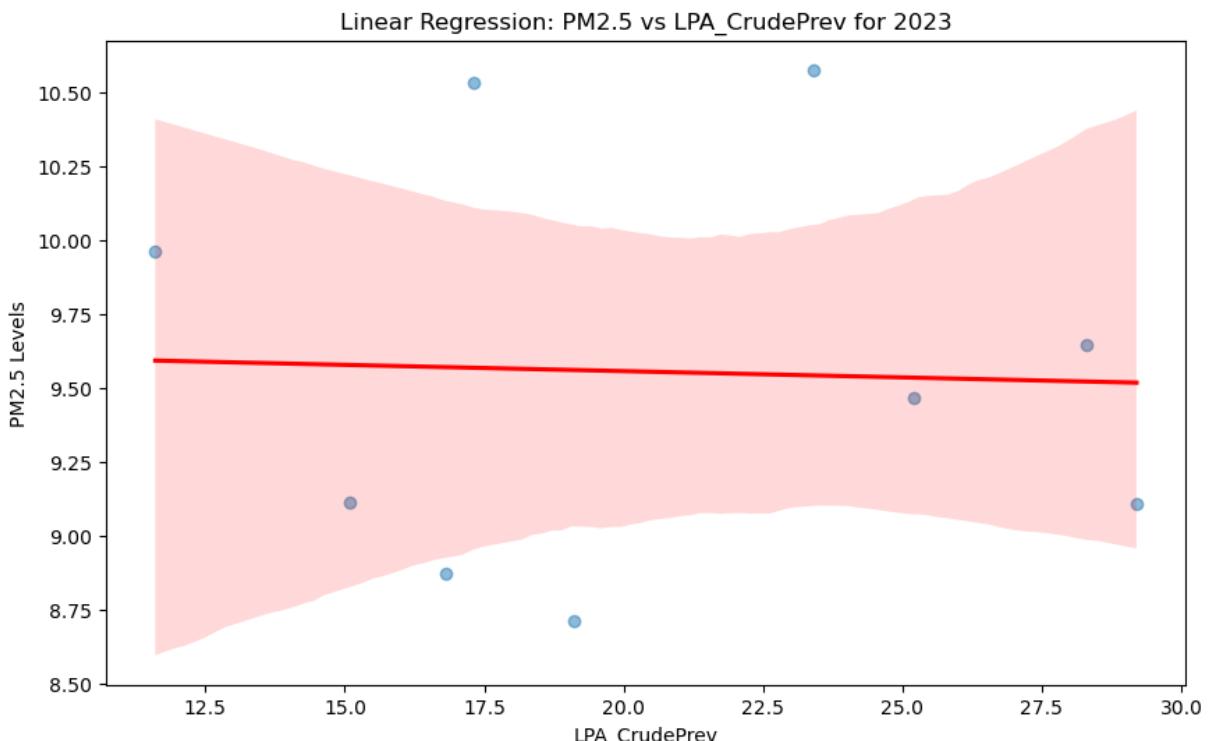
Linear Regression: PM2.5 vs DIABETES_CrudePrev for 2023



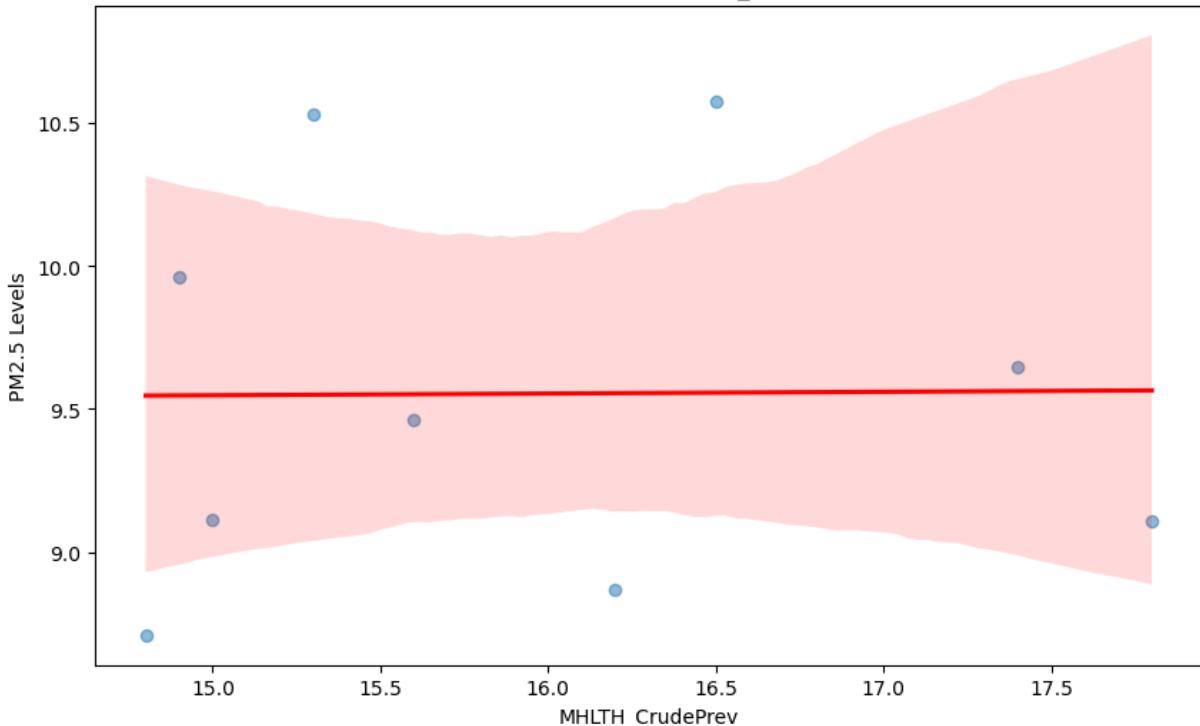
Linear Regression: PM2.5 vs GHLTH_CrudePrev for 2023



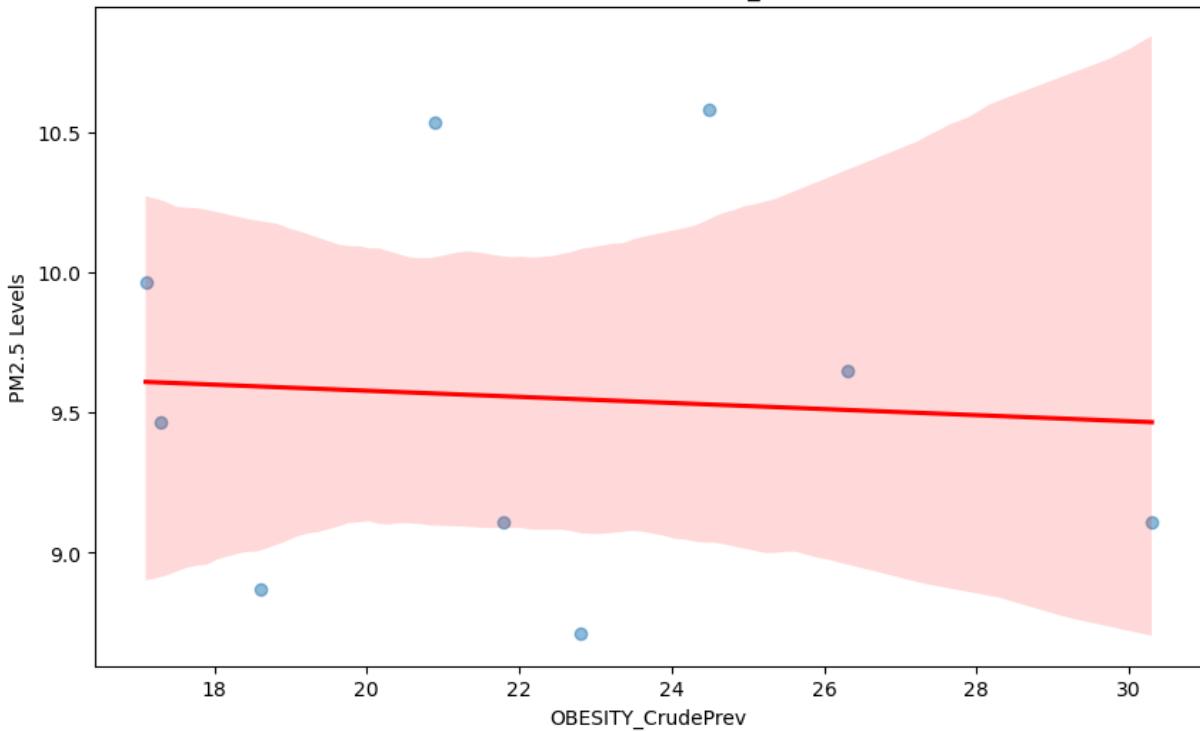


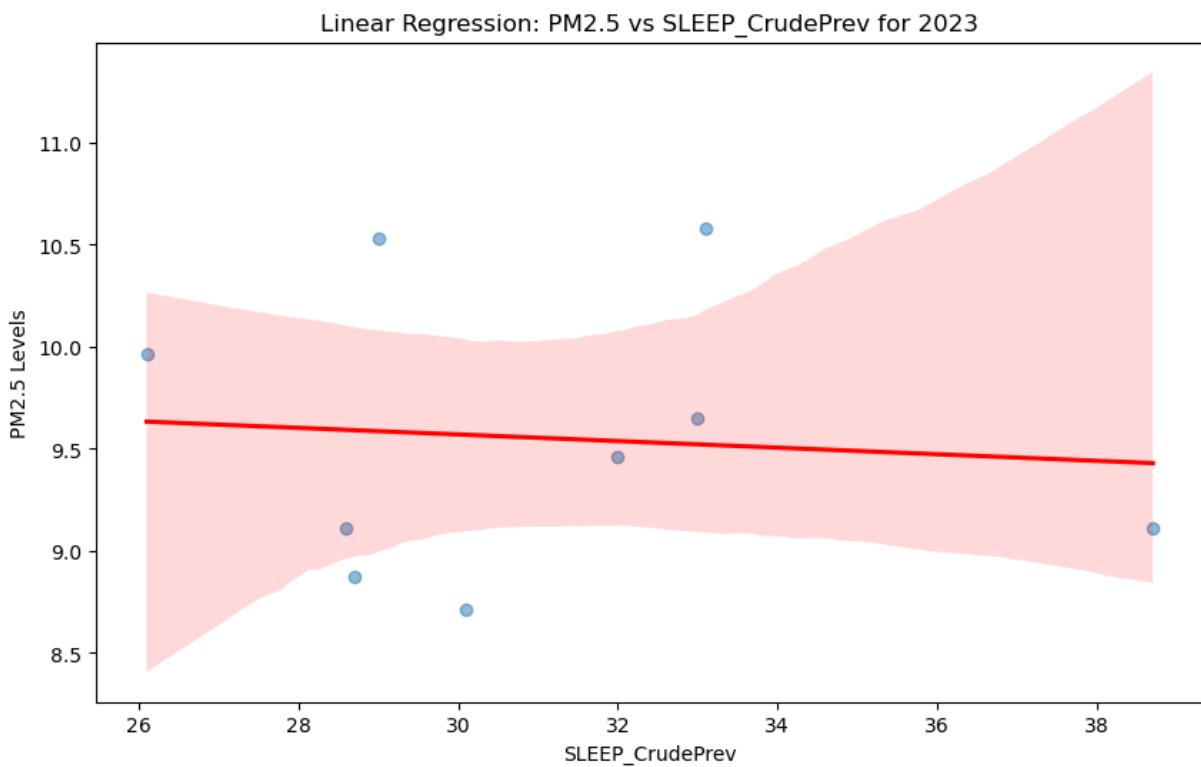
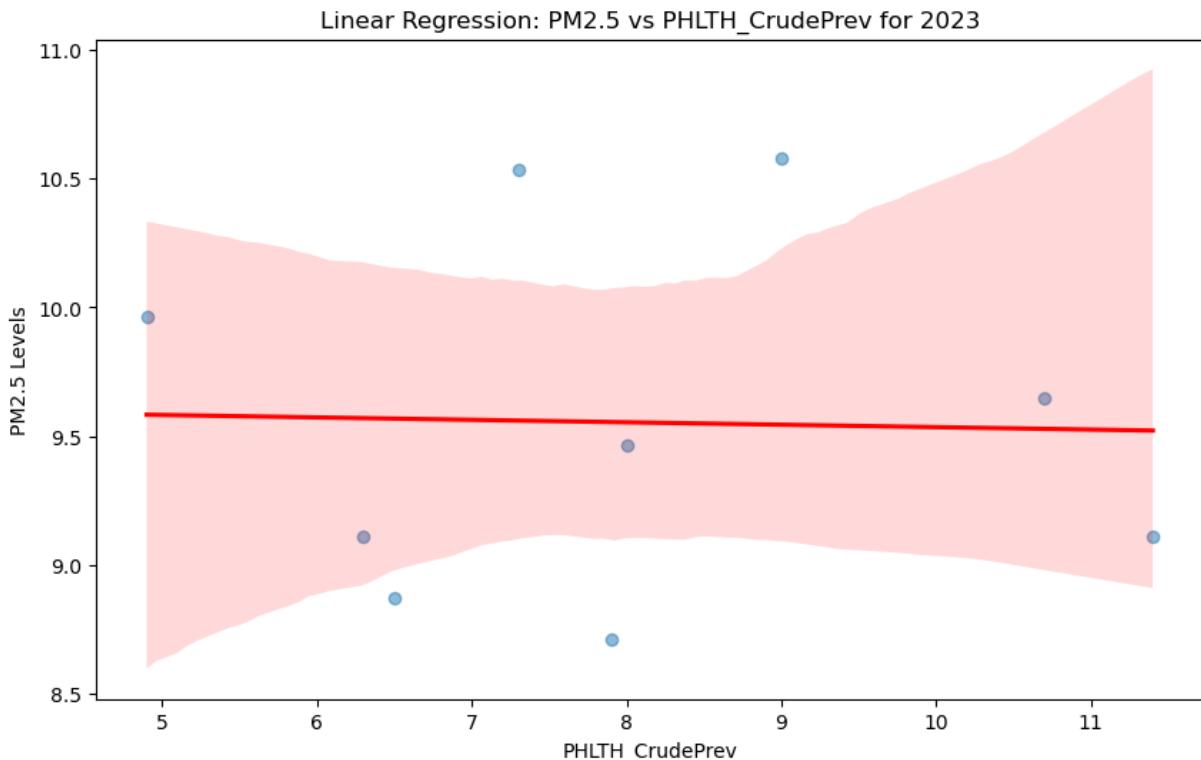


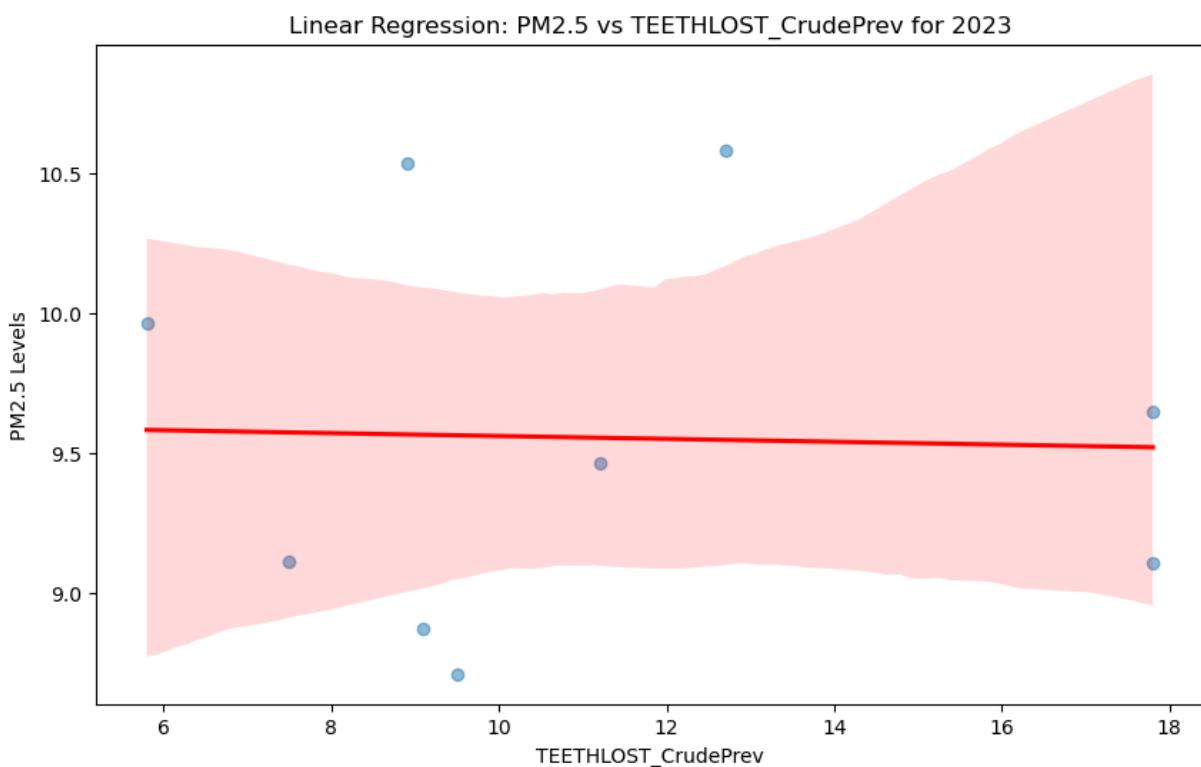
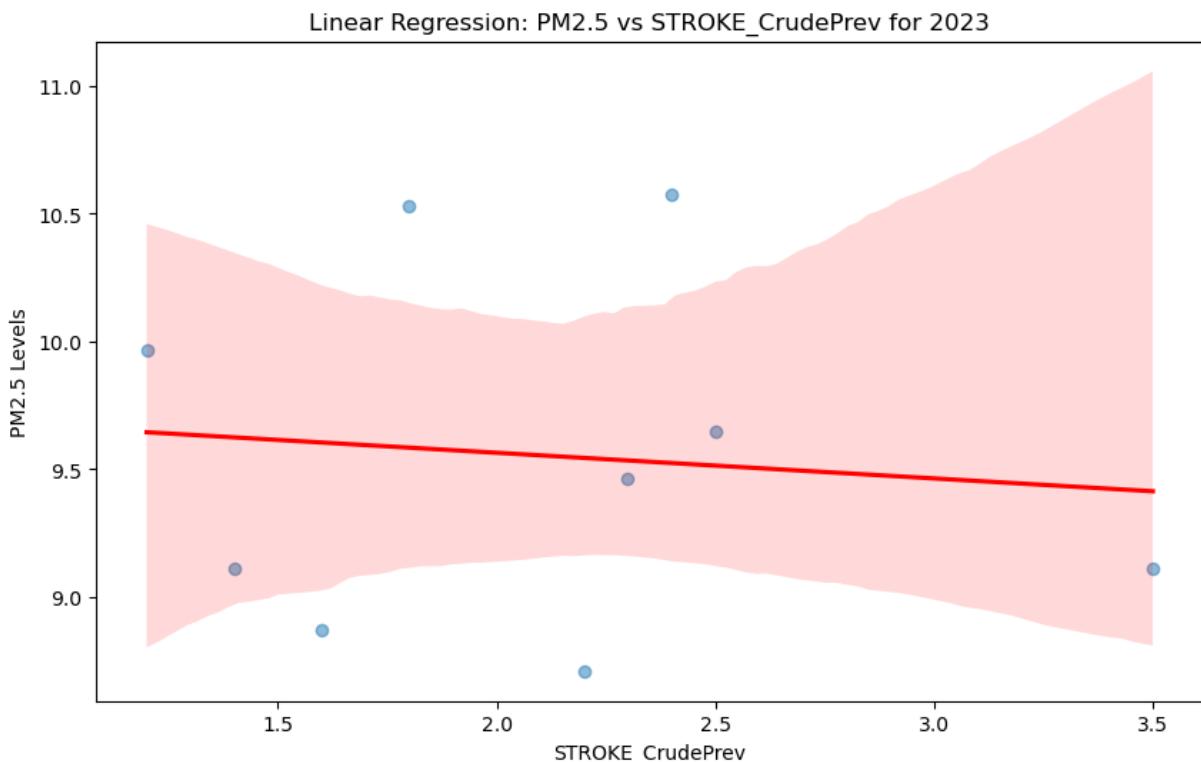
Linear Regression: PM2.5 vs MHLTH_CrudePrev for 2023

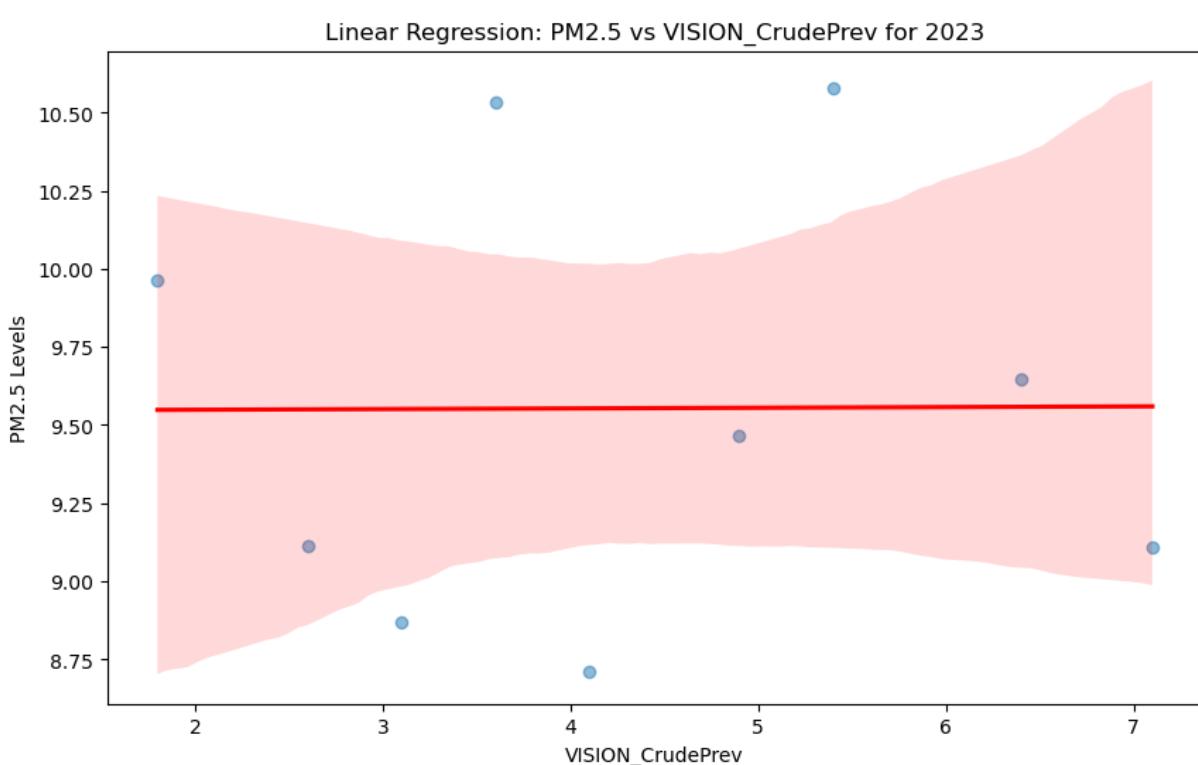
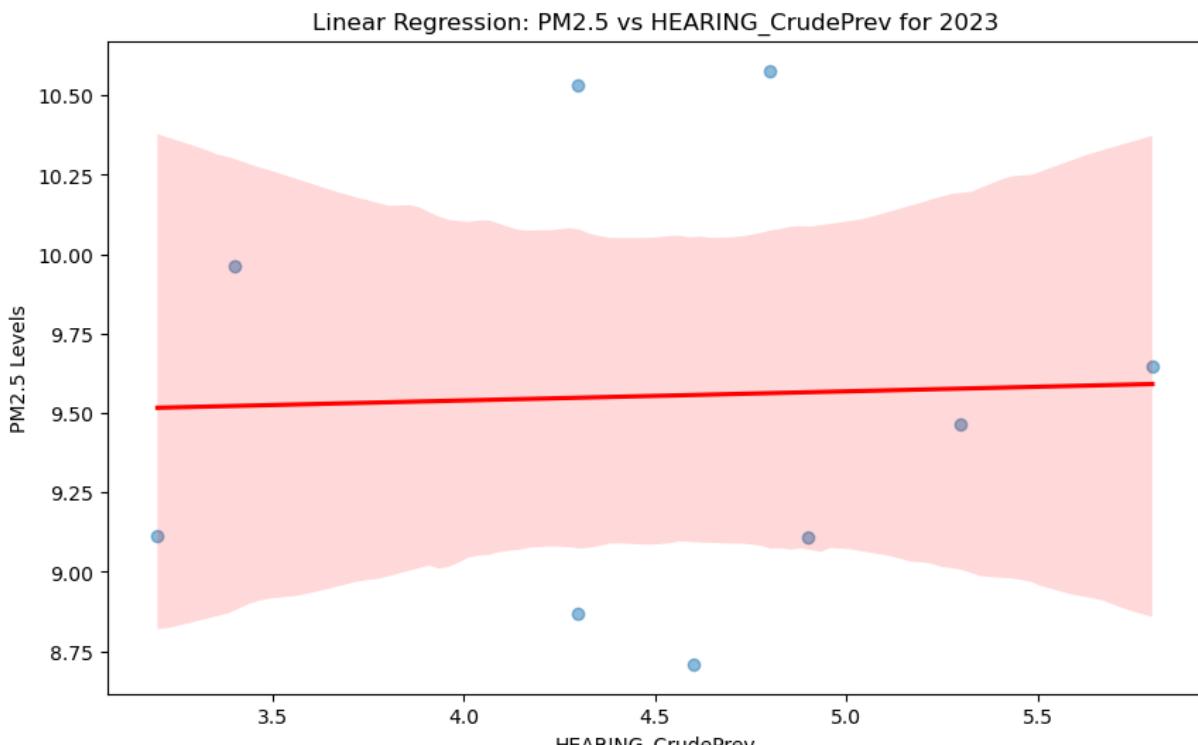


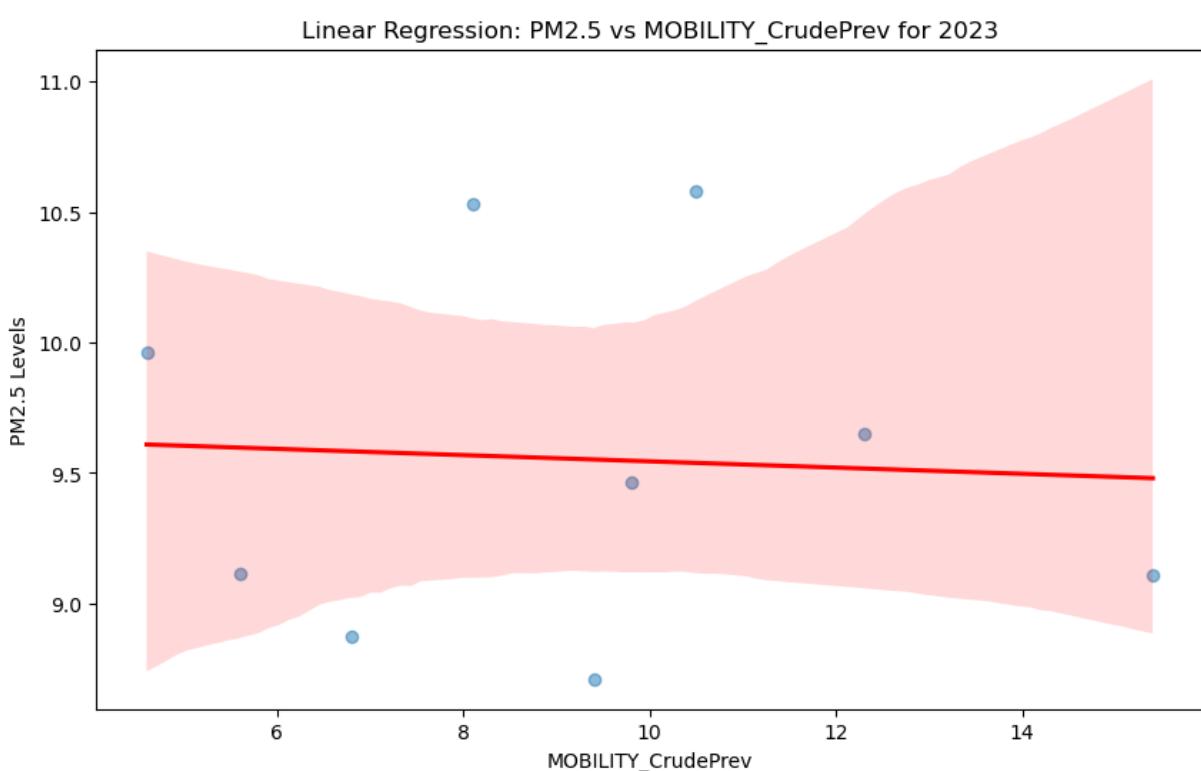
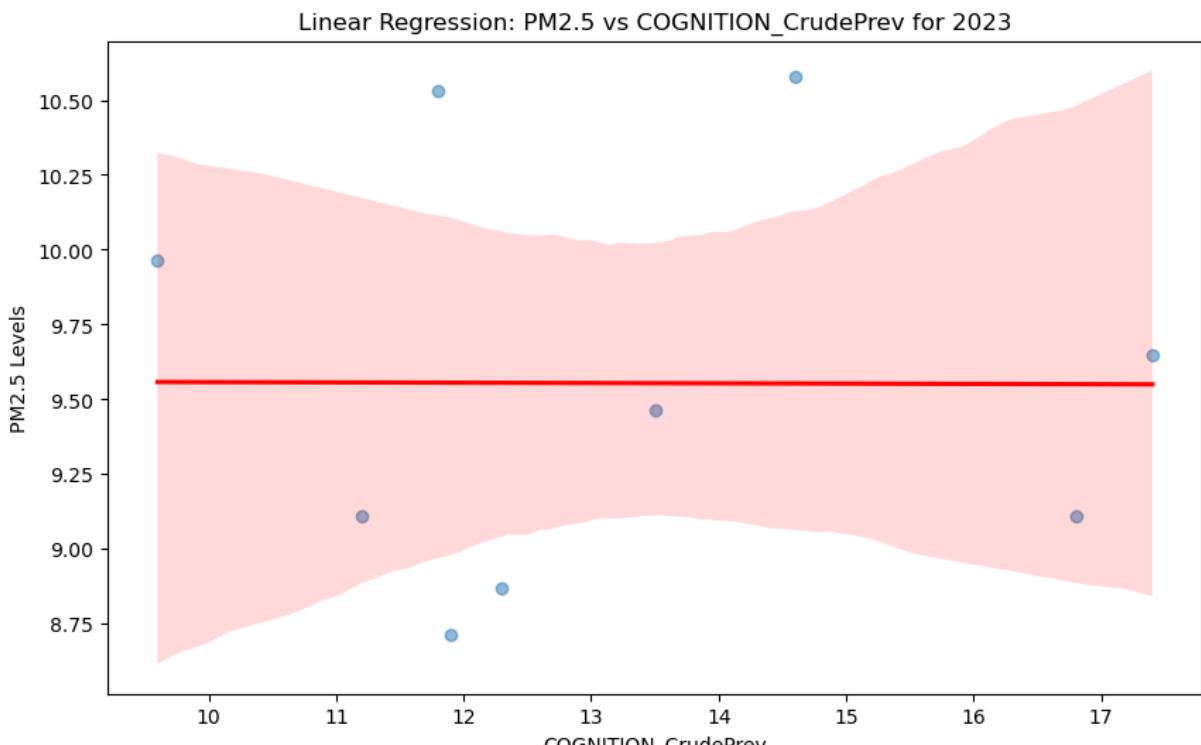
Linear Regression: PM2.5 vs OBESITY_CrudePrev for 2023



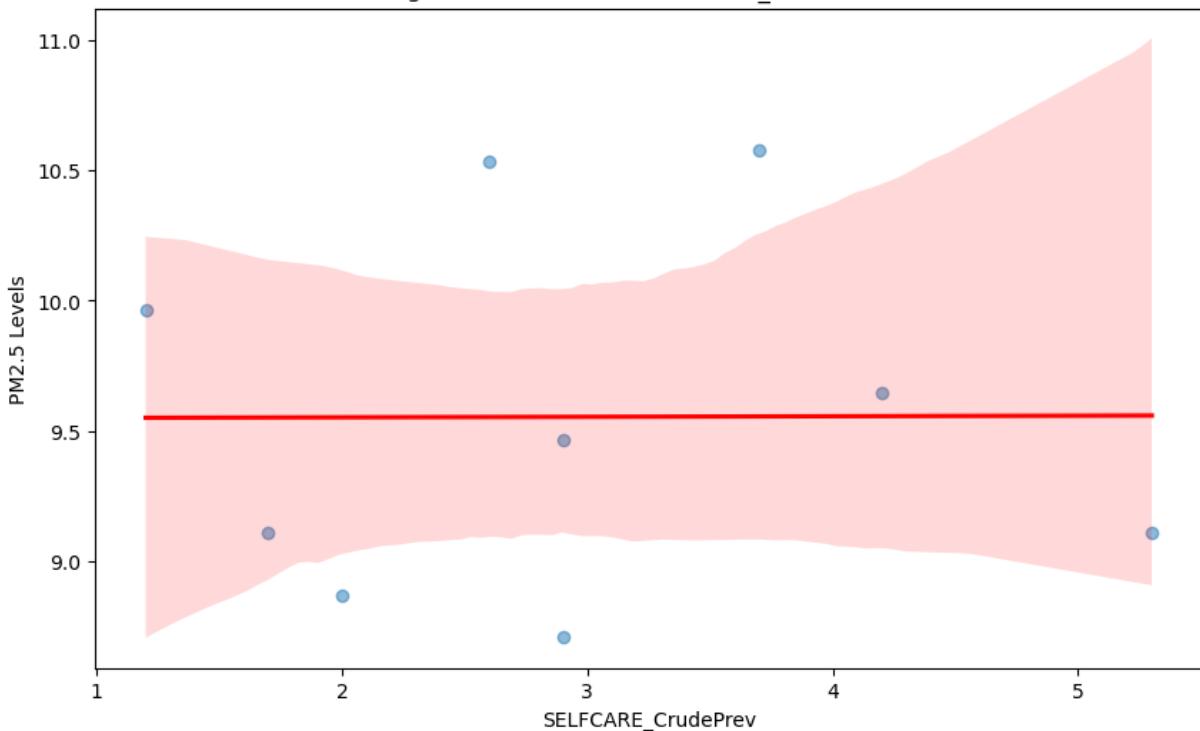




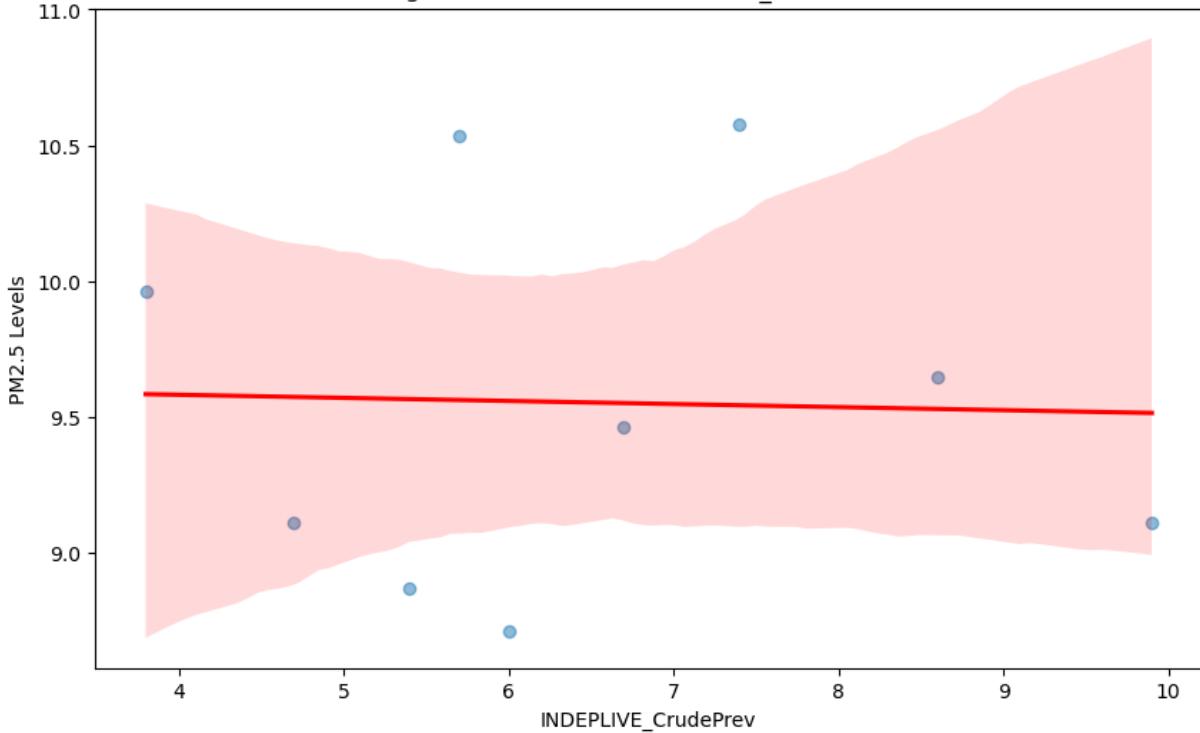


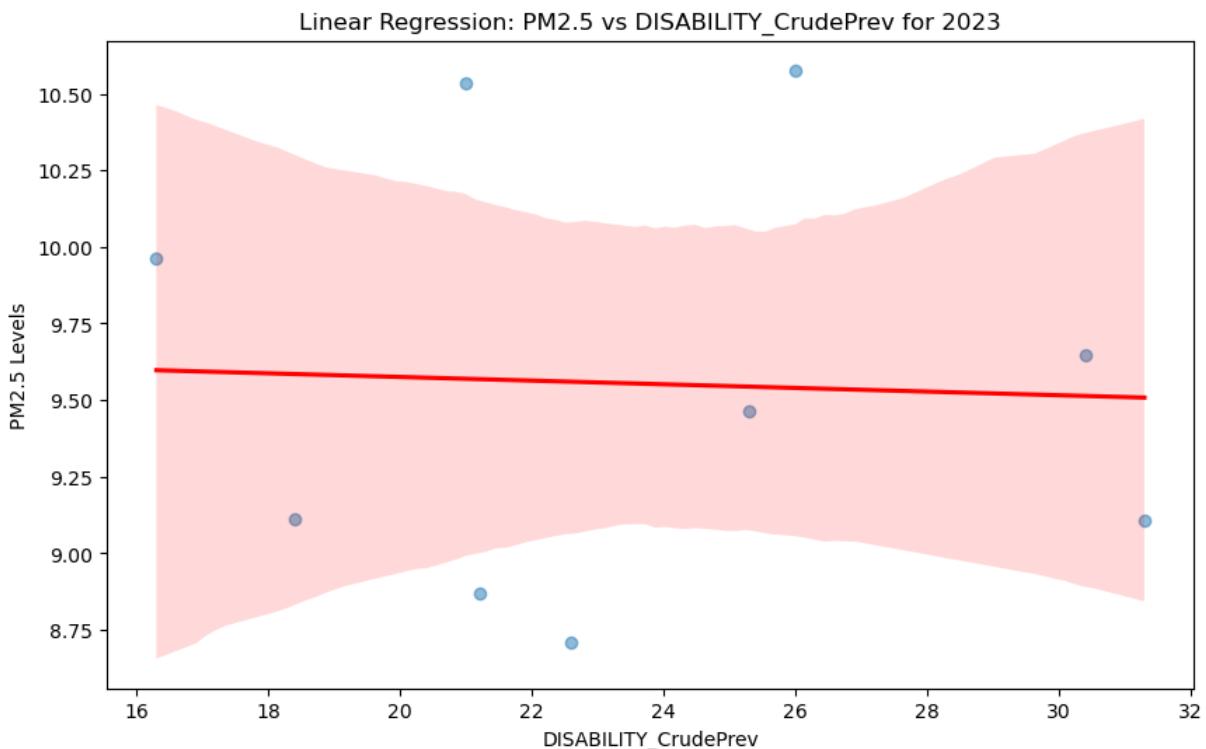


Linear Regression: PM2.5 vs SELFCARE_CrudePrev for 2023

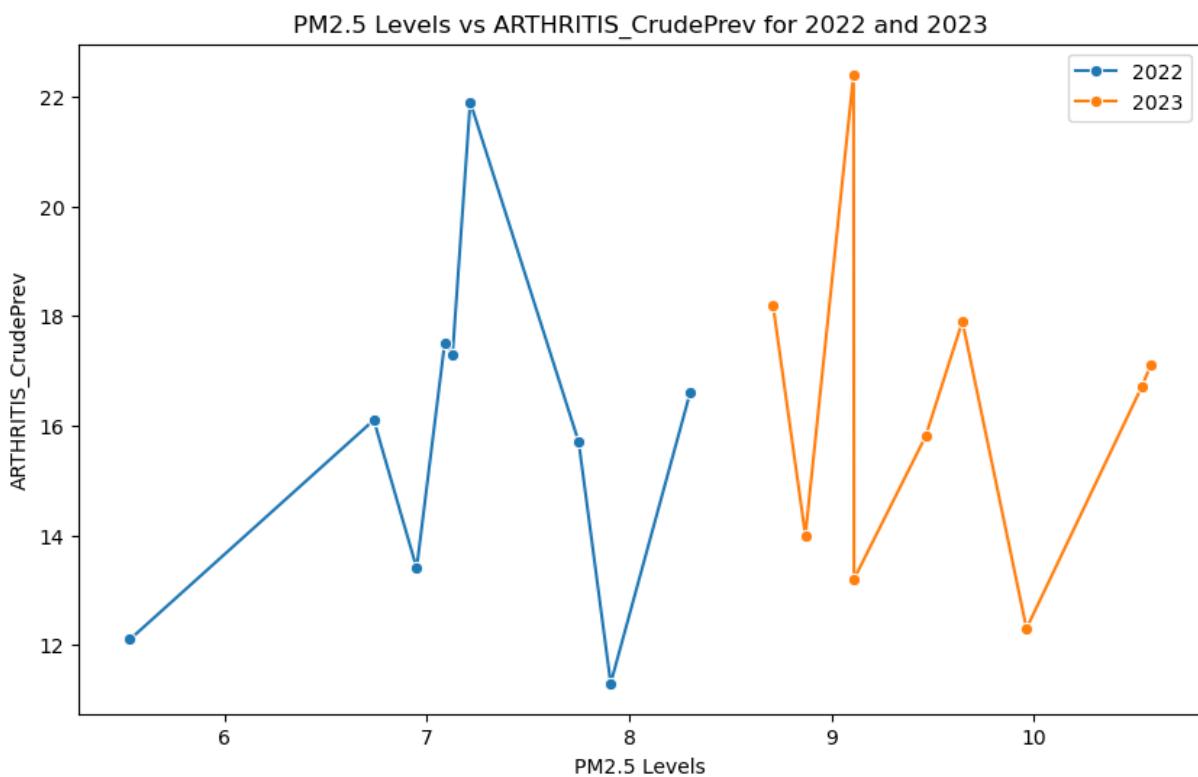
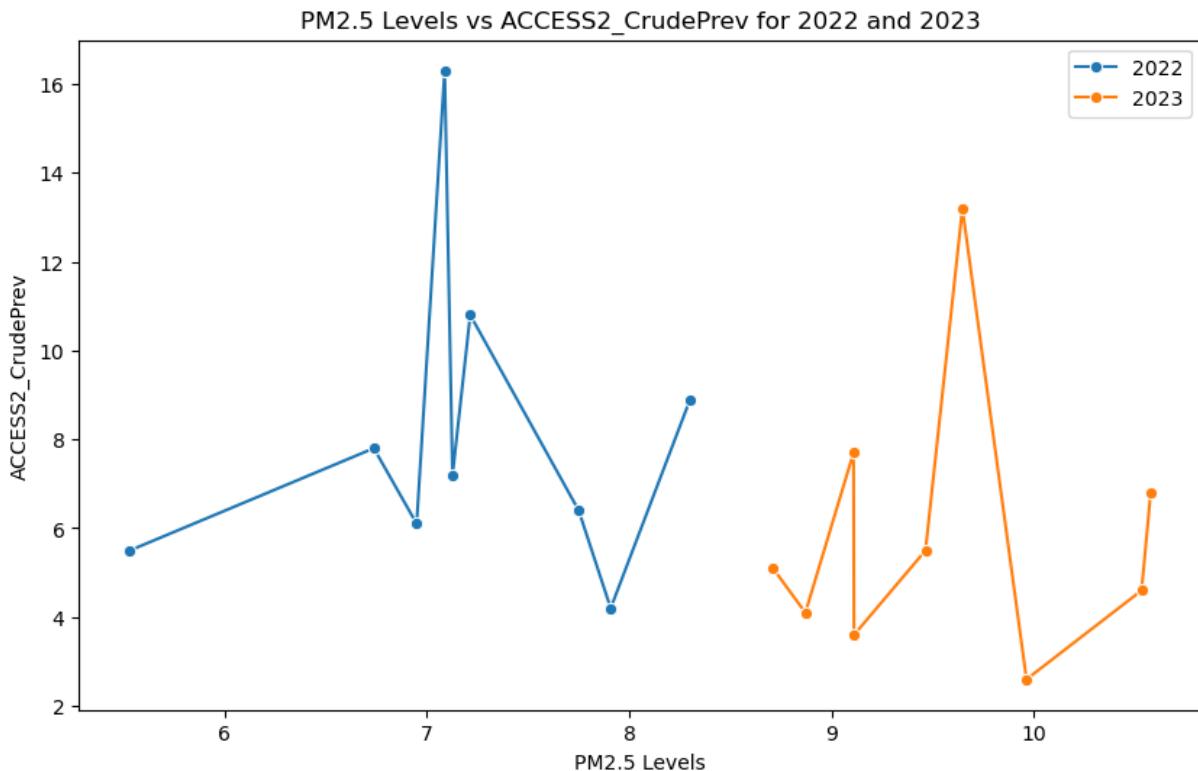


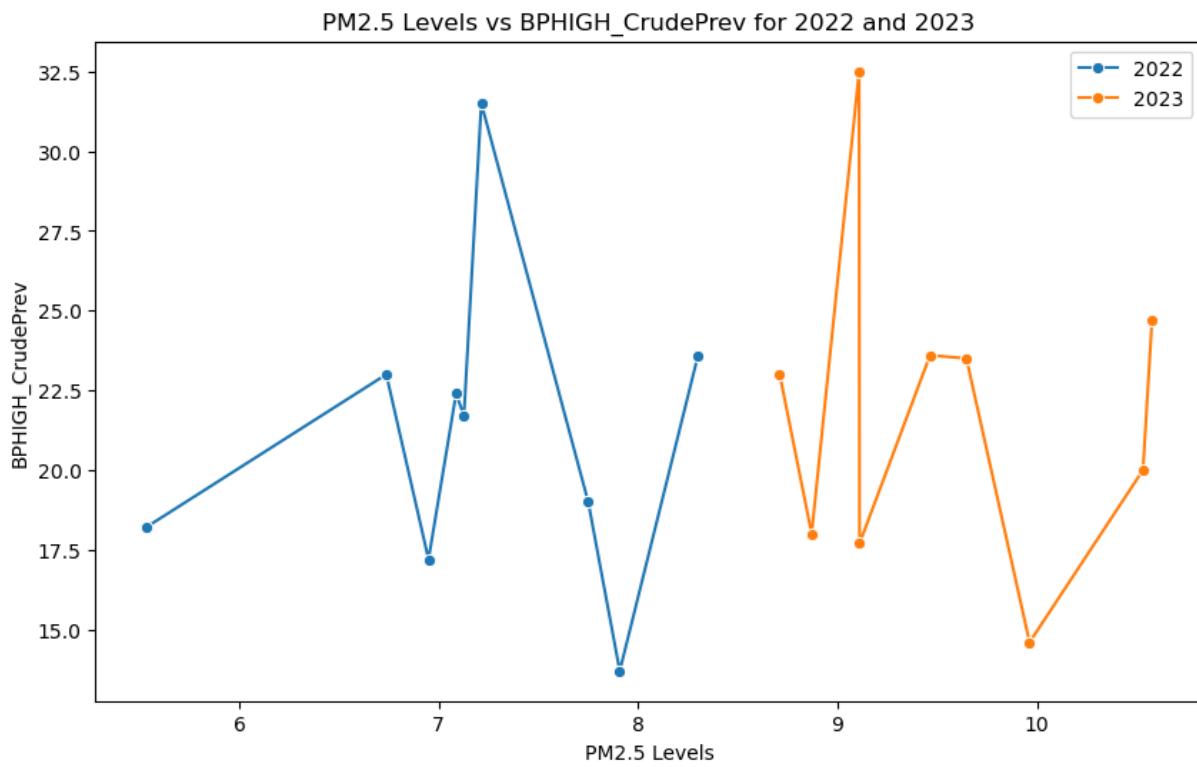
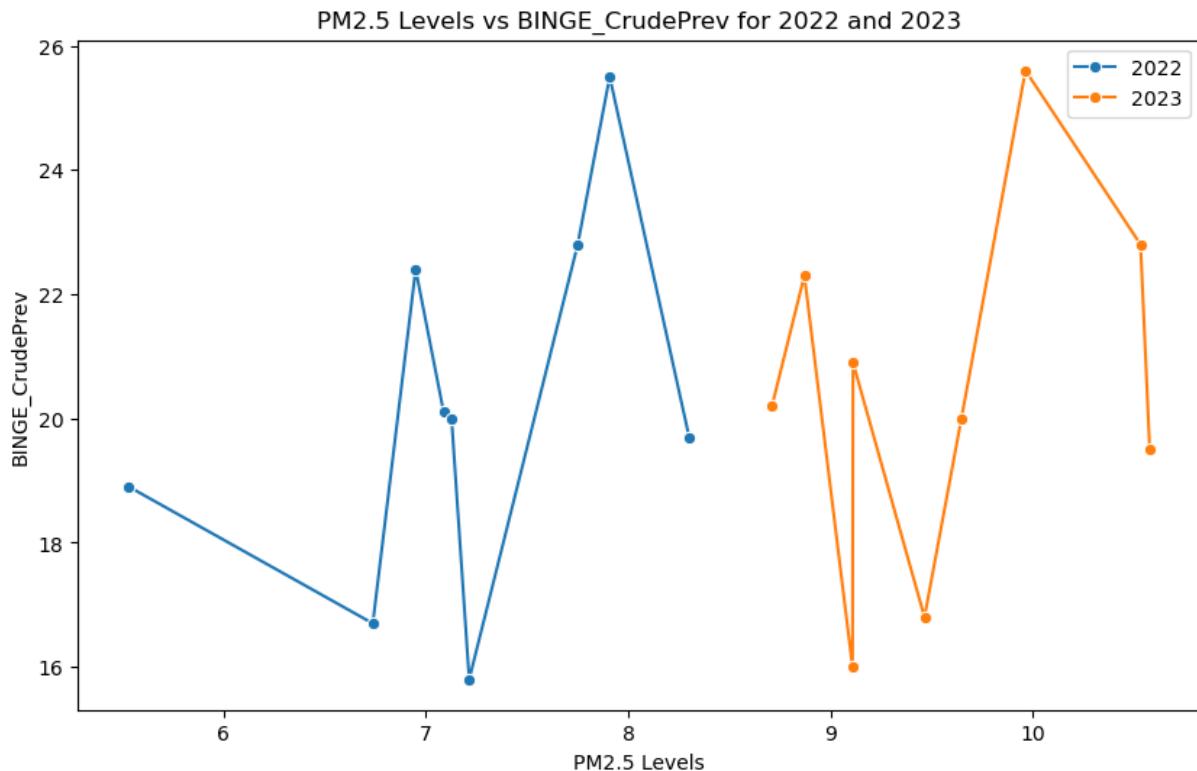
Linear Regression: PM2.5 vs INDEPLIVE_CrudePrev for 2023

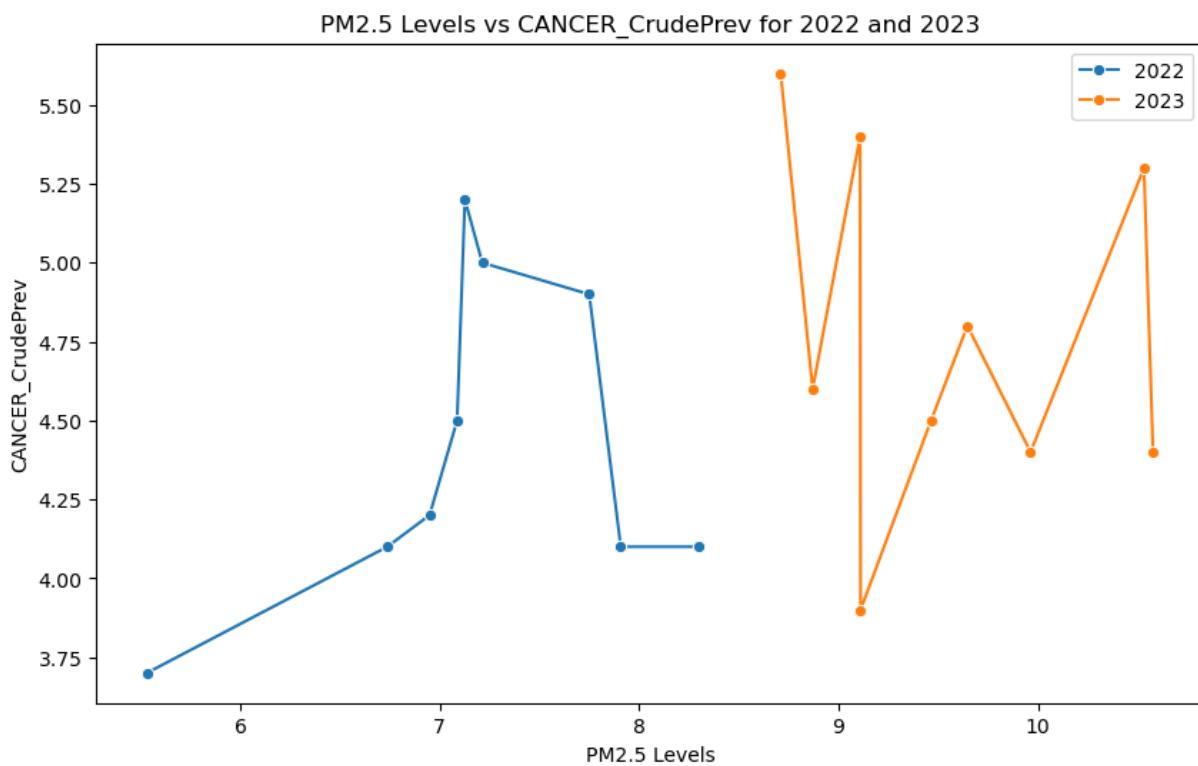
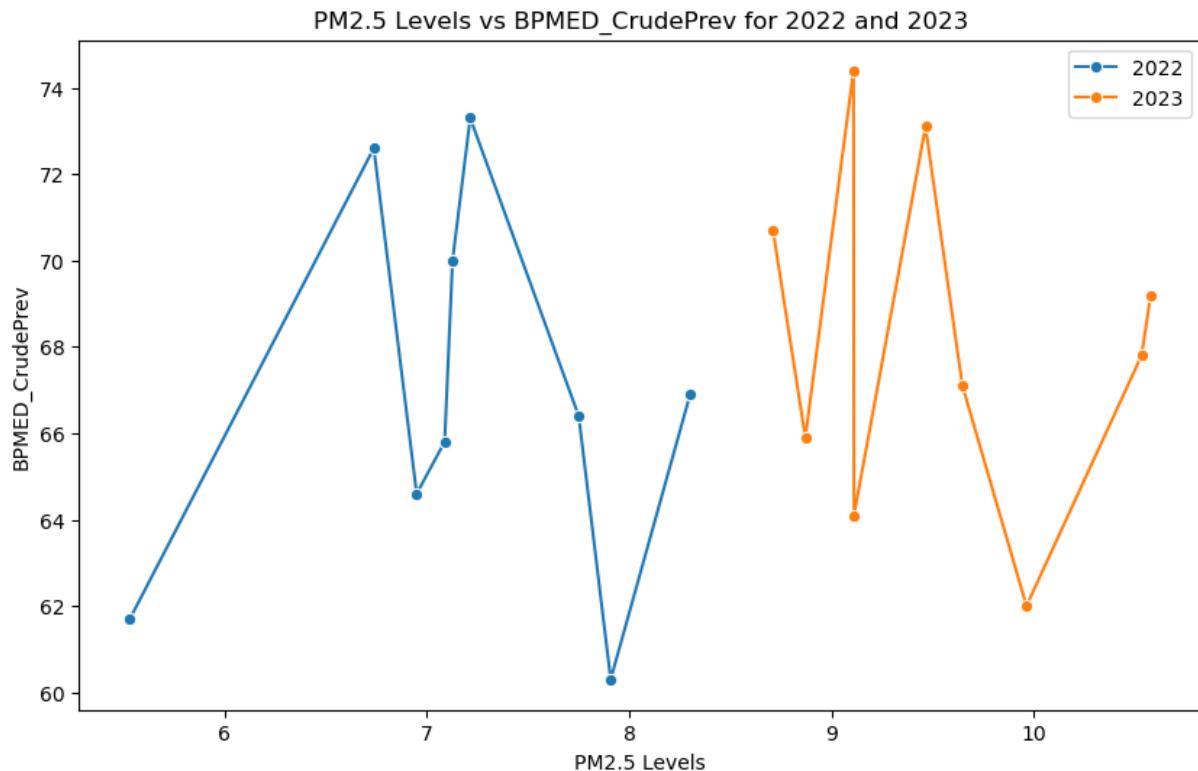


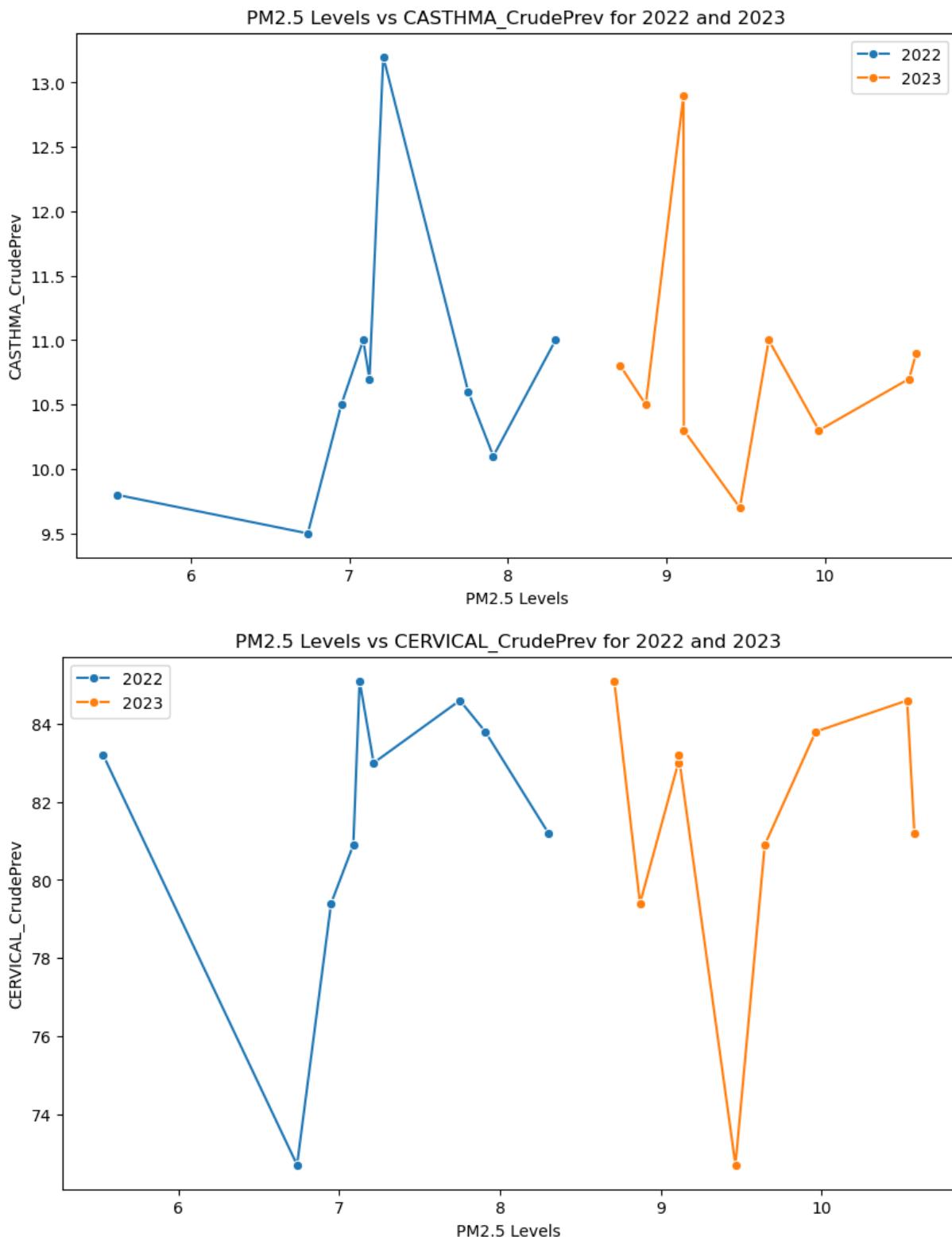


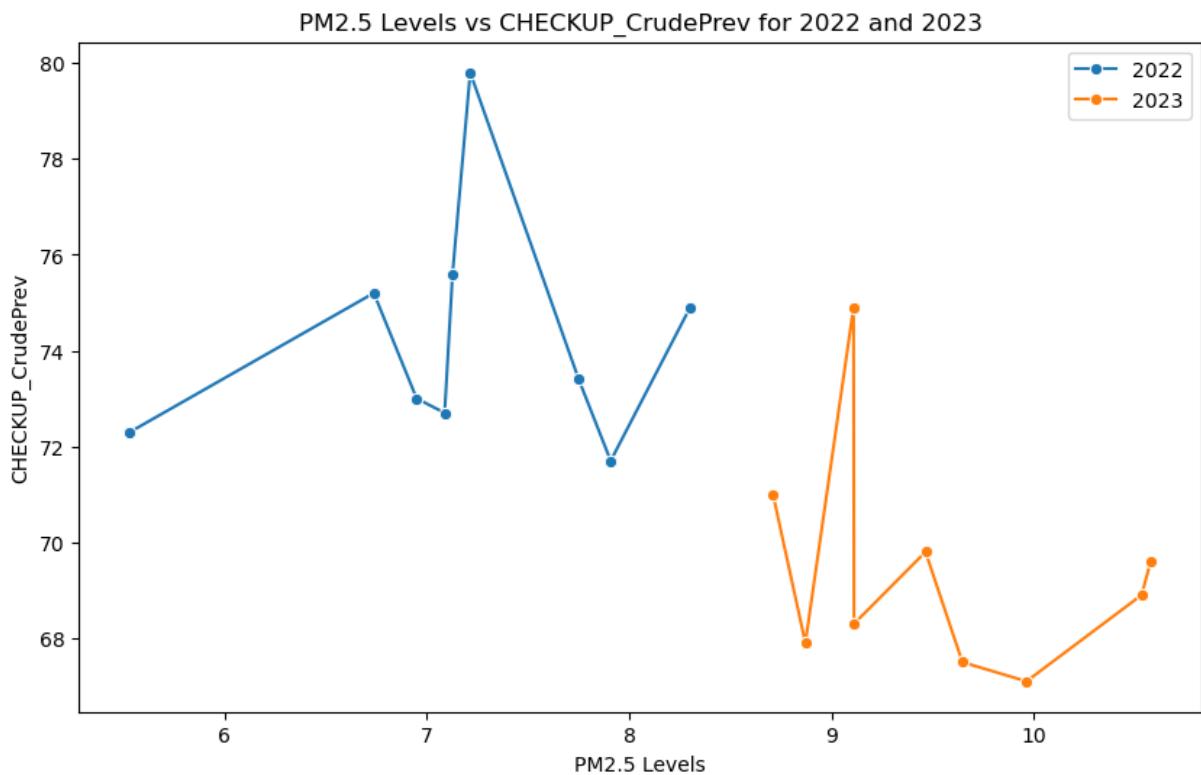
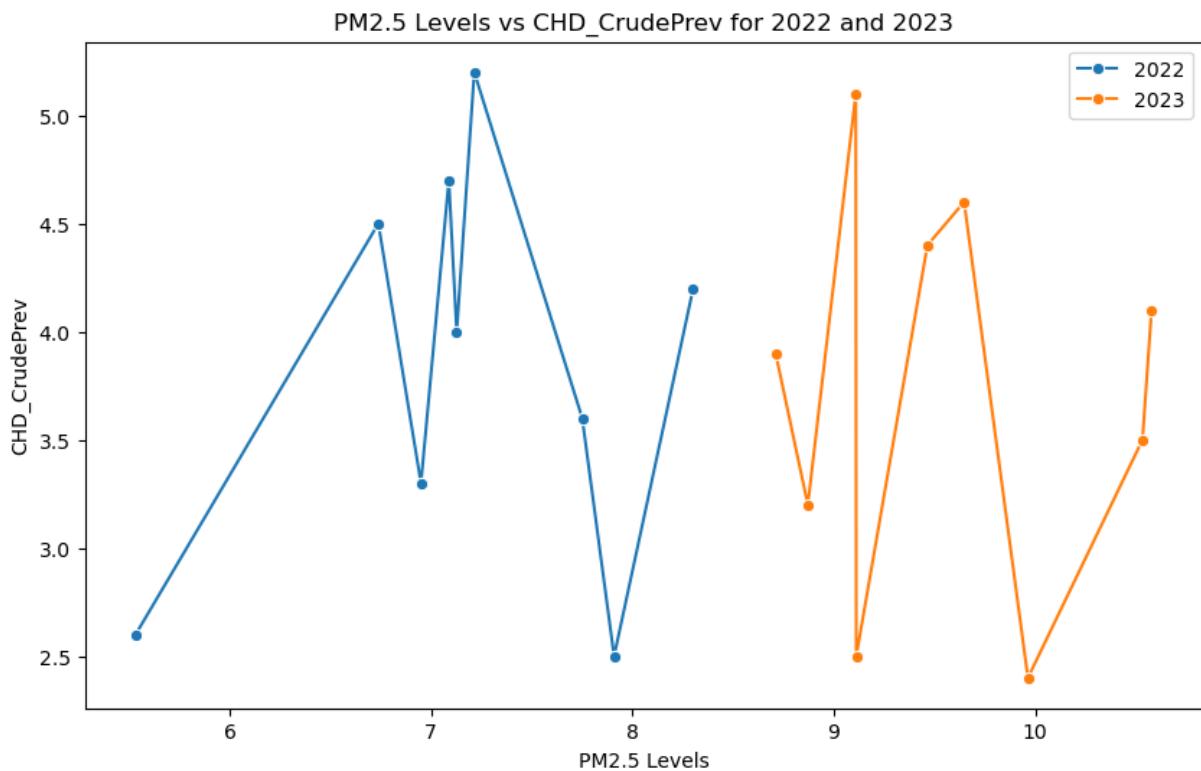
```
In [11]: health_data_2023_renamed = health_data_2023.rename(columns=lambda x: x.replace('Crude', 'CrudePrev'))  
  
# Merging the datasets again with corrected column names  
merged_health_data_corrected = pd.merge(health_data_2022, health_data_2023_renamed)  
  
# Plotting line graphs for each crude rate against PM2.5 level for 2022 and 2023  
for column in crude_columns_2022:  
    if 'Crude' in column:  
        plt.figure(figsize=(10, 6))  
  
        # Plotting for 2022  
        sns.lineplot(x=merged_health_data_corrected['PM2.5_2022'], y=merged_health_data_corrected[column])  
  
        # Plotting for 2023  
        sns.lineplot(x=merged_health_data_corrected['PM2.5_2023'], y=merged_health_data_corrected[column])  
  
        plt.title(f'PM2.5 Levels vs {column} for 2022 and 2023')  
        plt.xlabel('PM2.5 Levels')  
        plt.ylabel(column)  
        plt.legend()  
        plt.show()
```

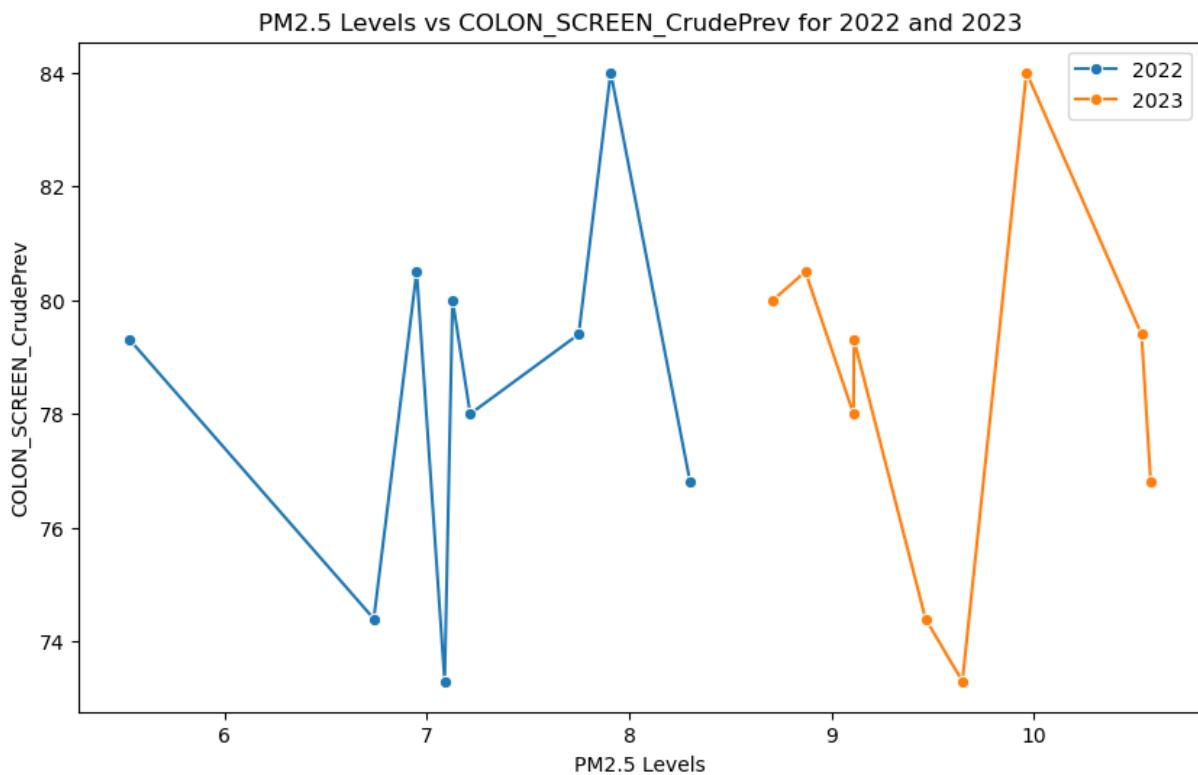
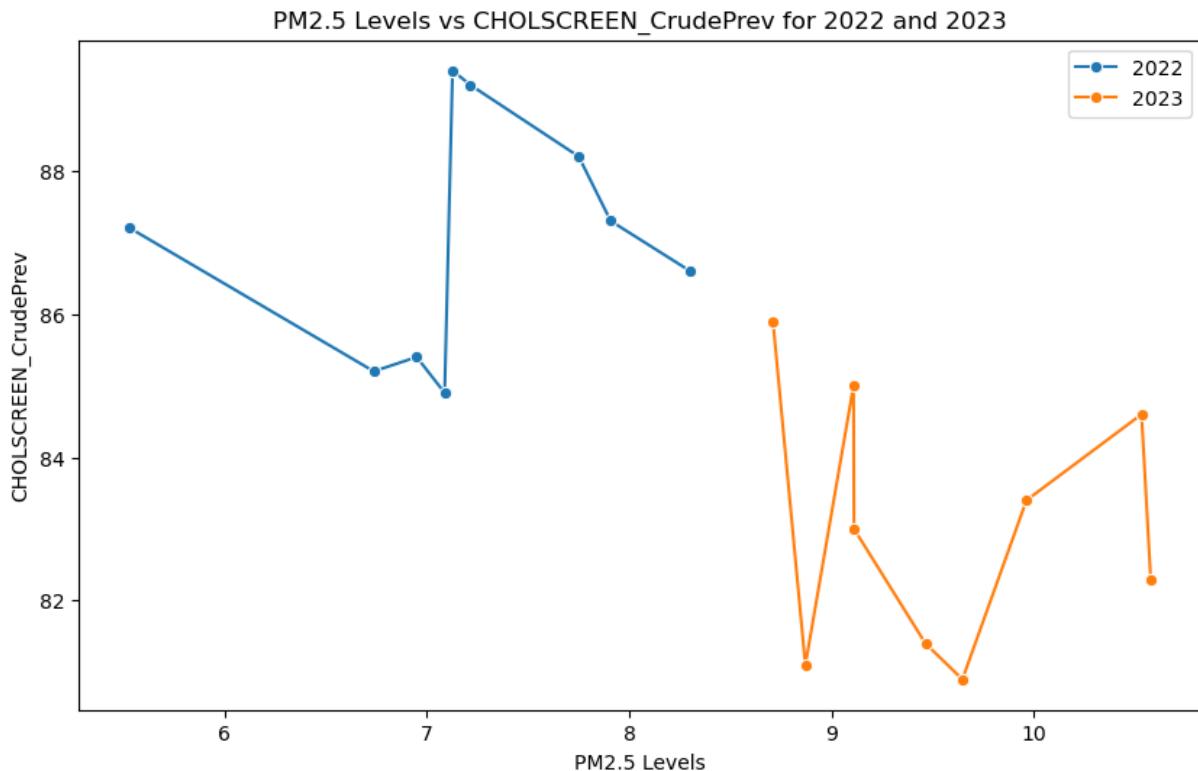




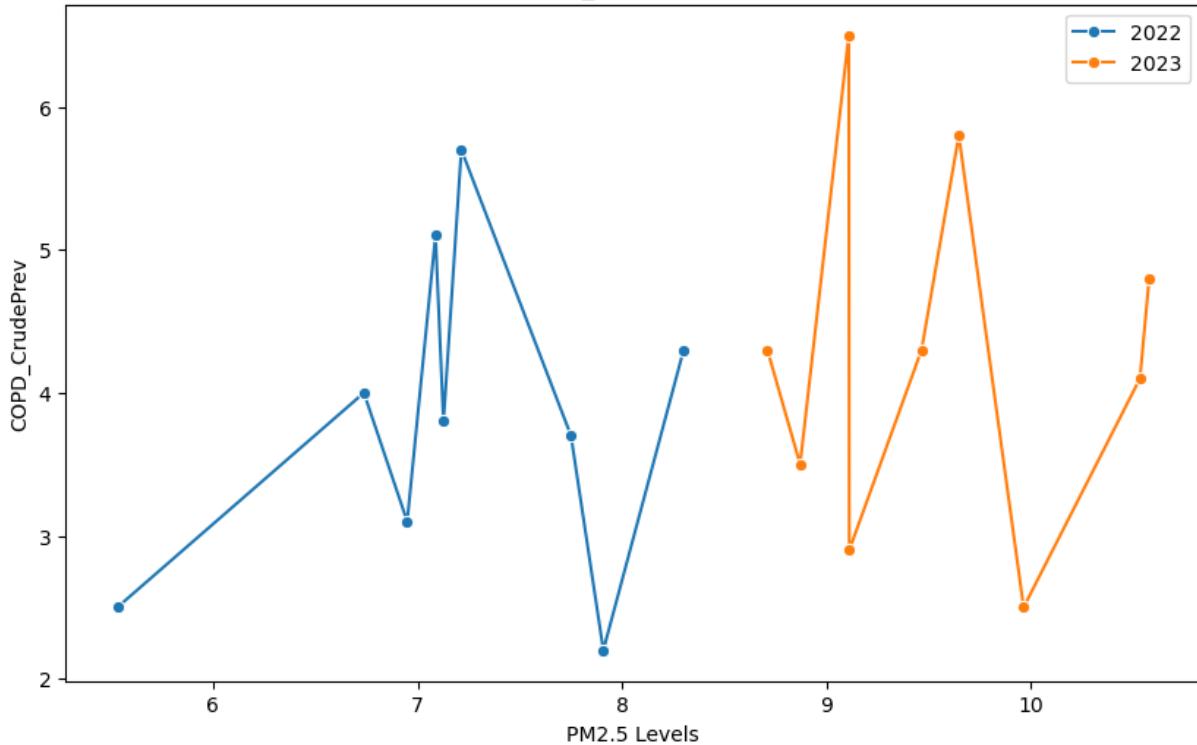




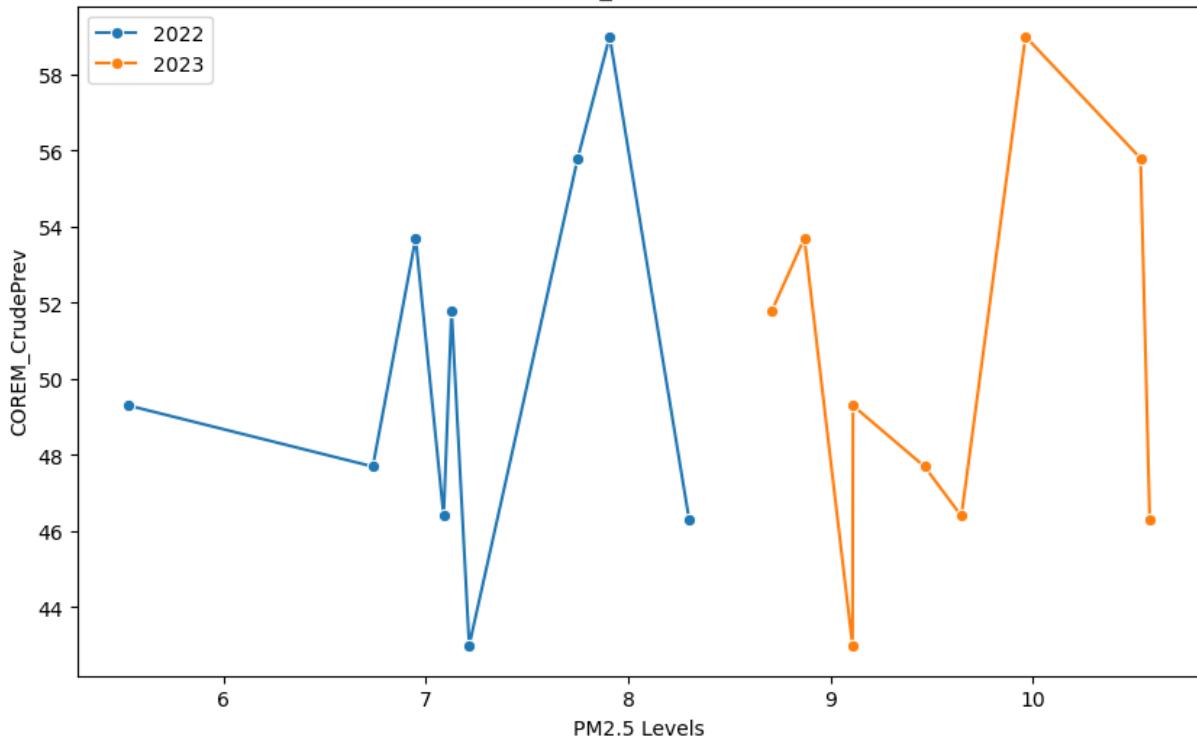




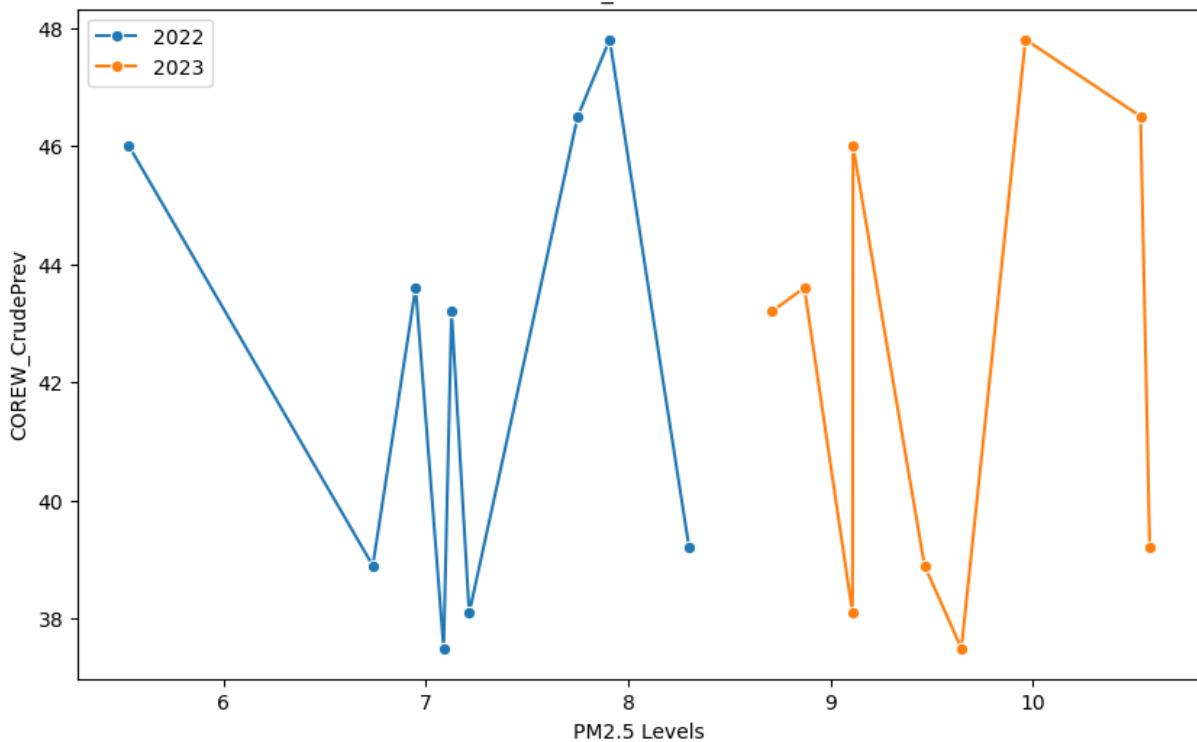
PM2.5 Levels vs COPD_CrudePrev for 2022 and 2023



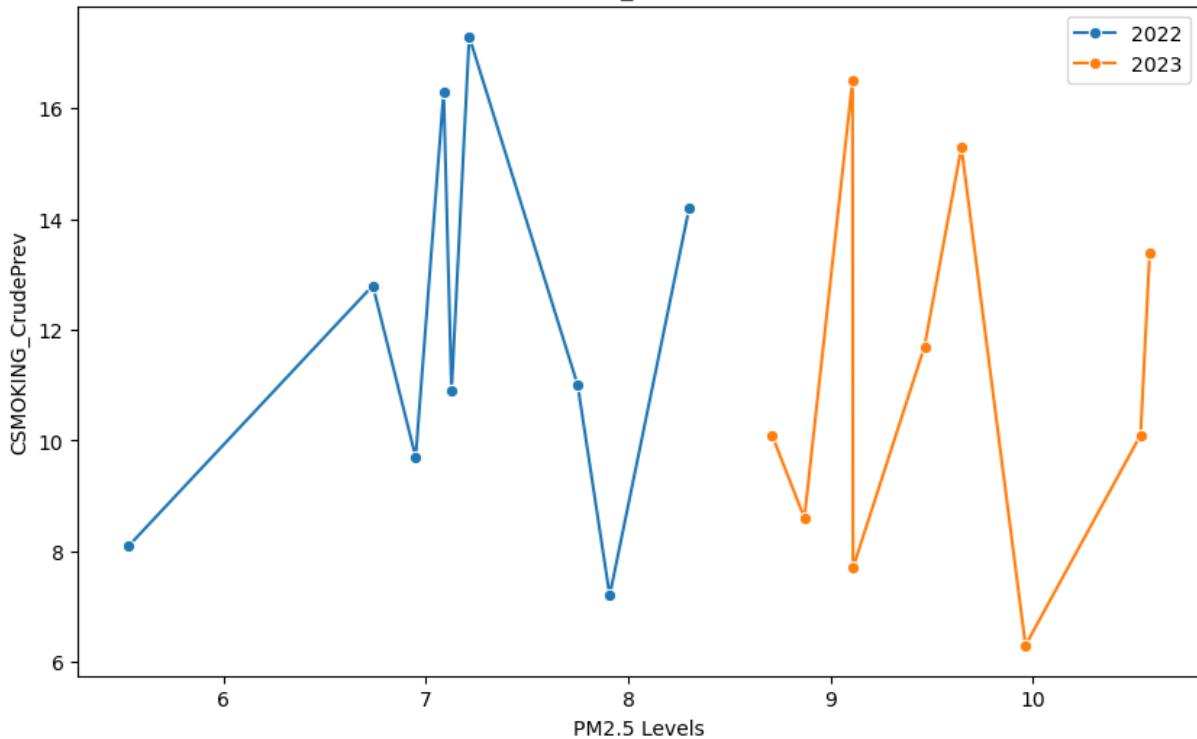
PM2.5 Levels vs COREM_CrudePrev for 2022 and 2023

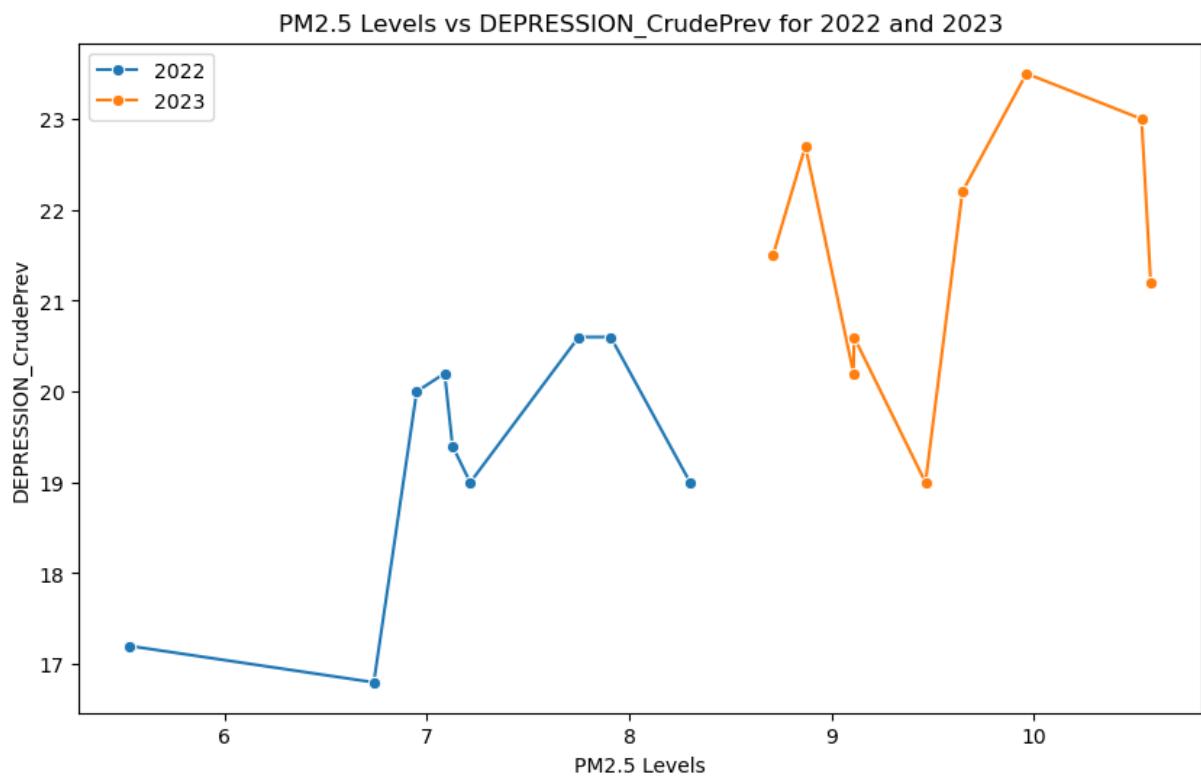
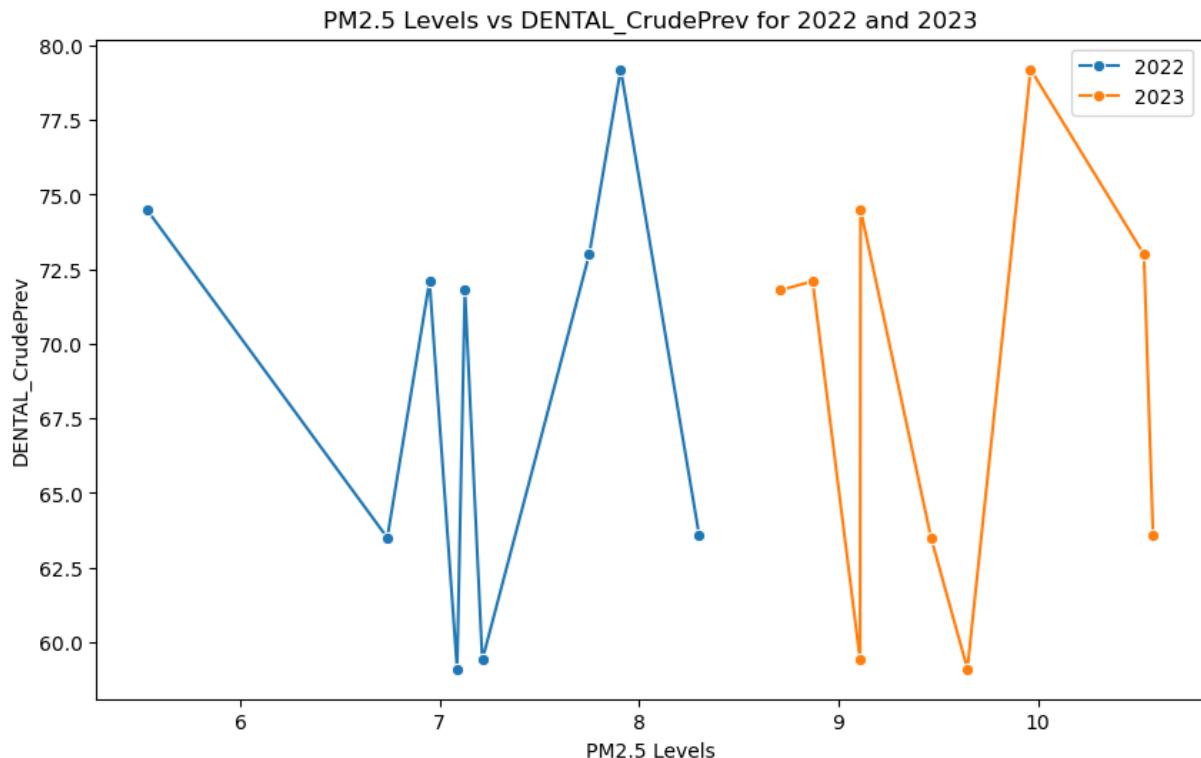


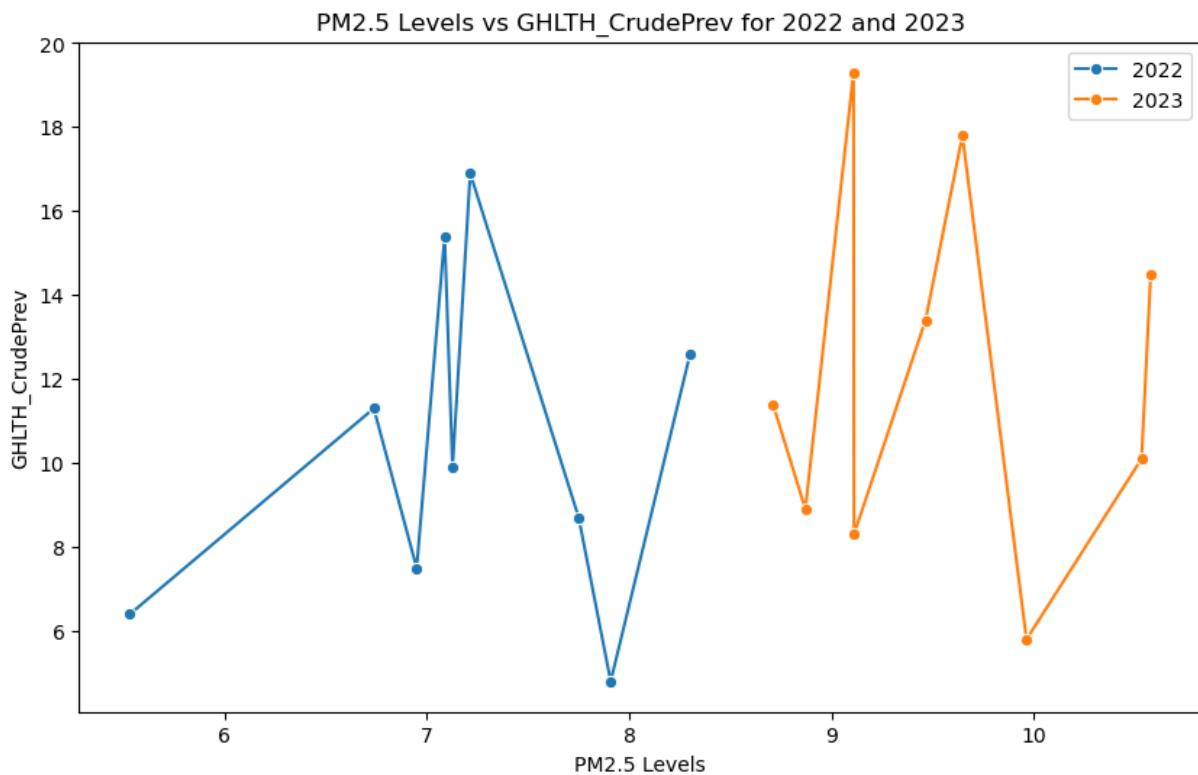
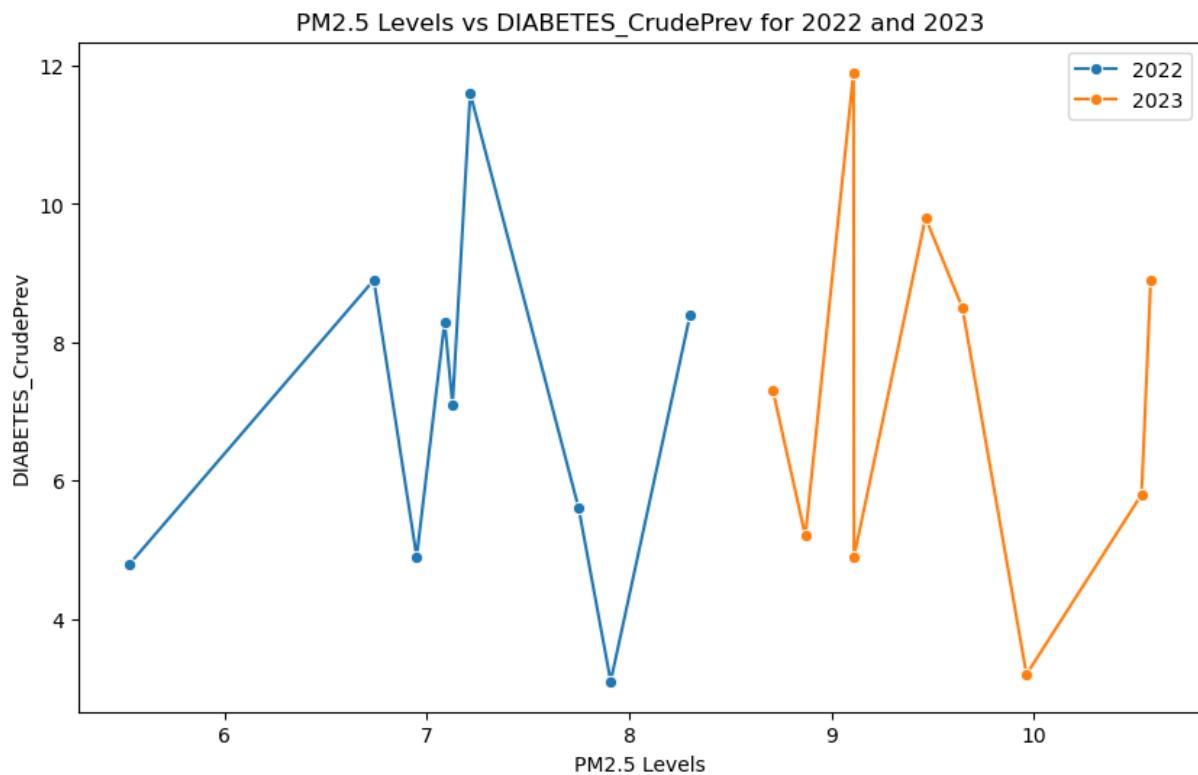
PM2.5 Levels vs COREW_CrudePrev for 2022 and 2023

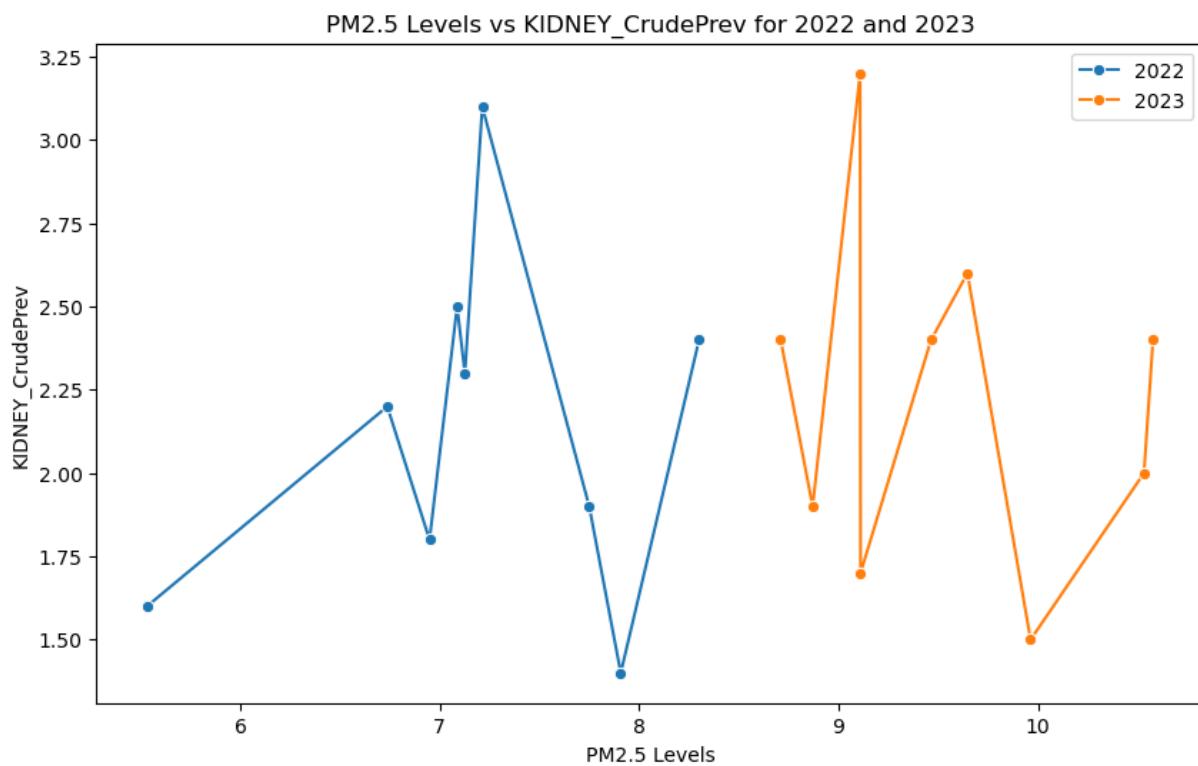
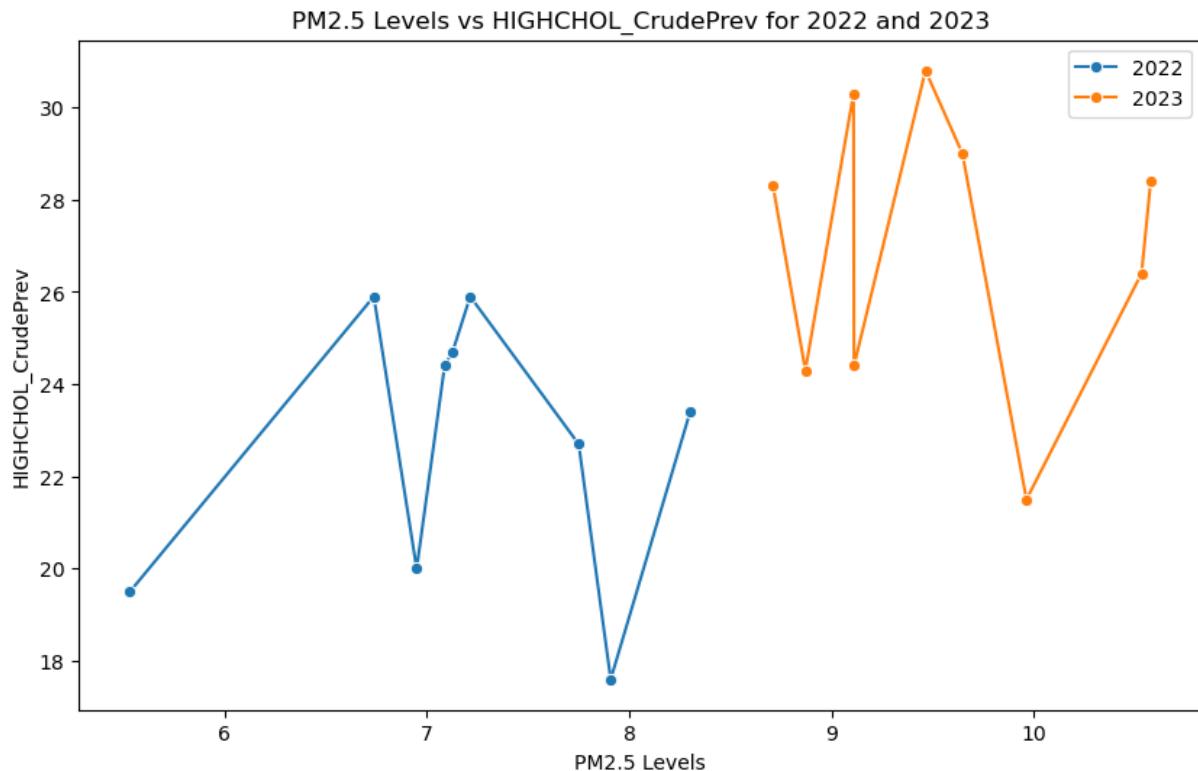


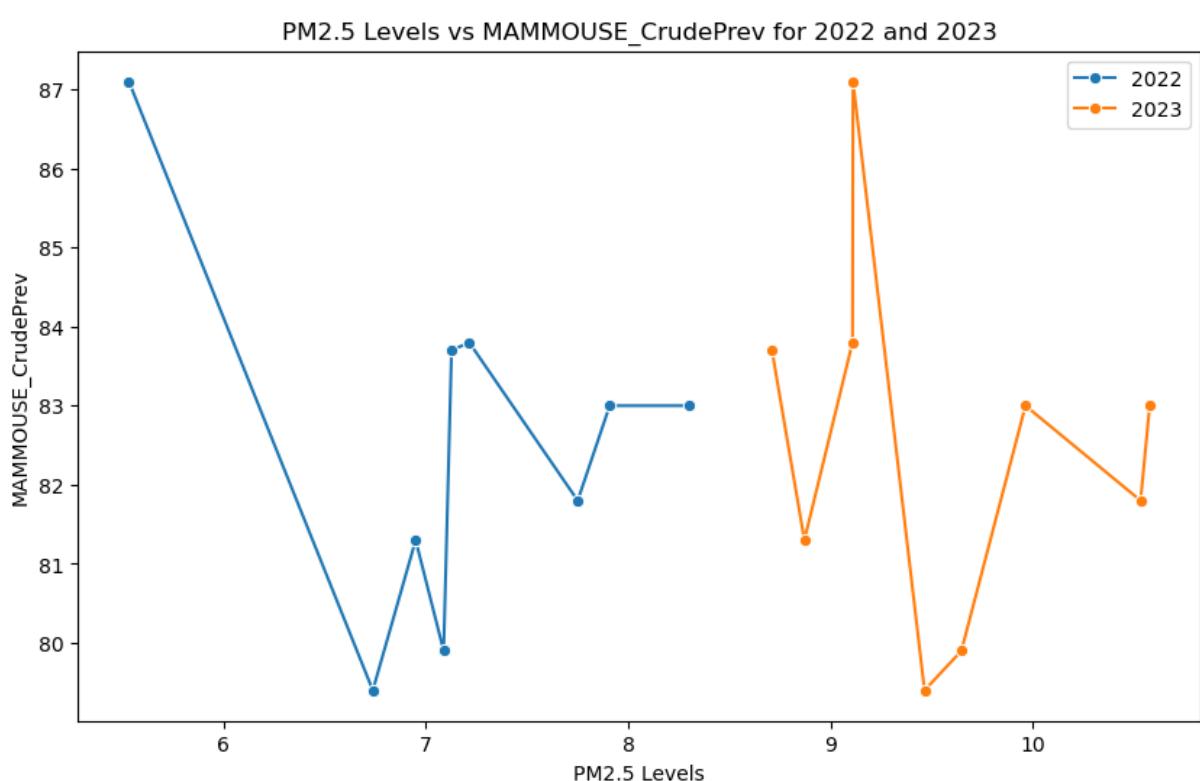
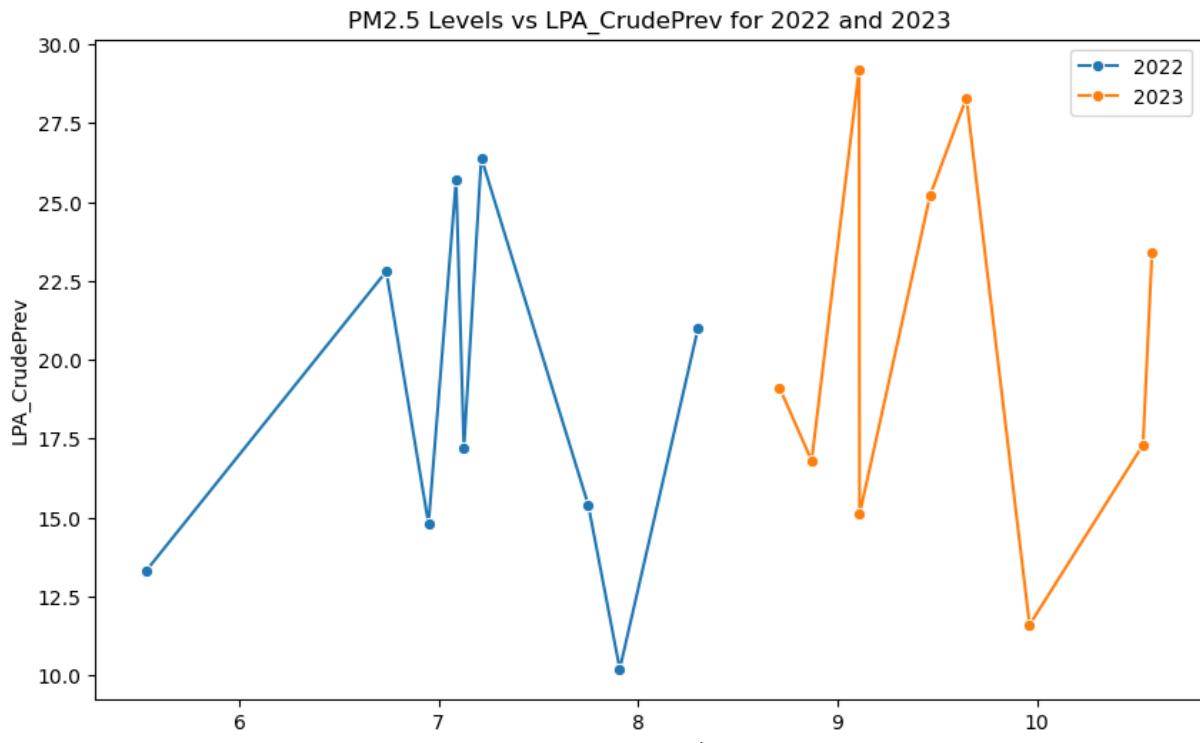
PM2.5 Levels vs CSMOKING_CrudePrev for 2022 and 2023

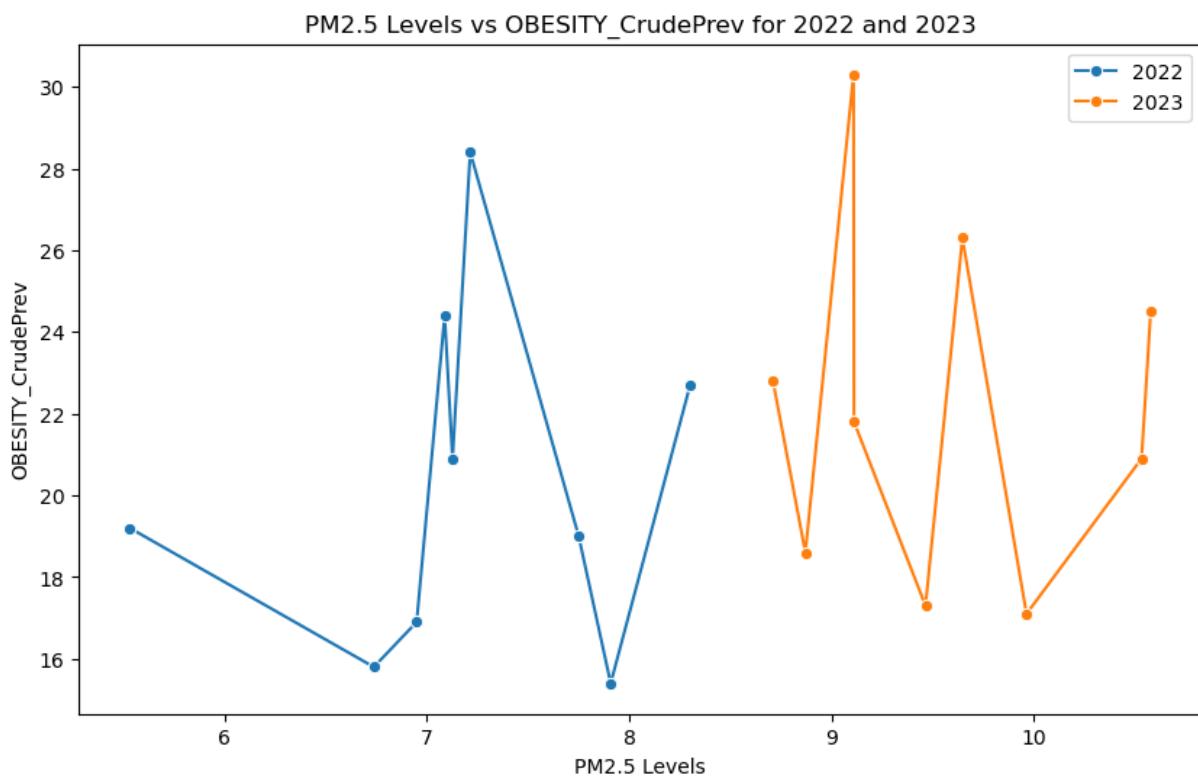
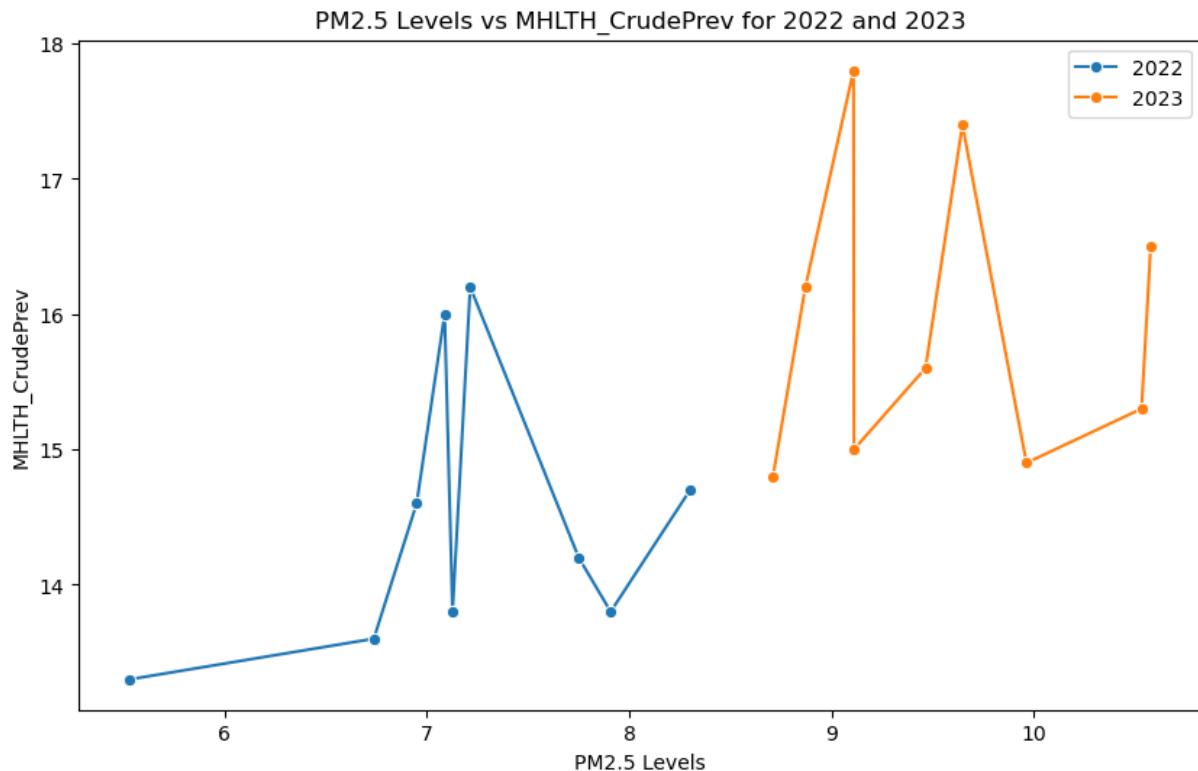


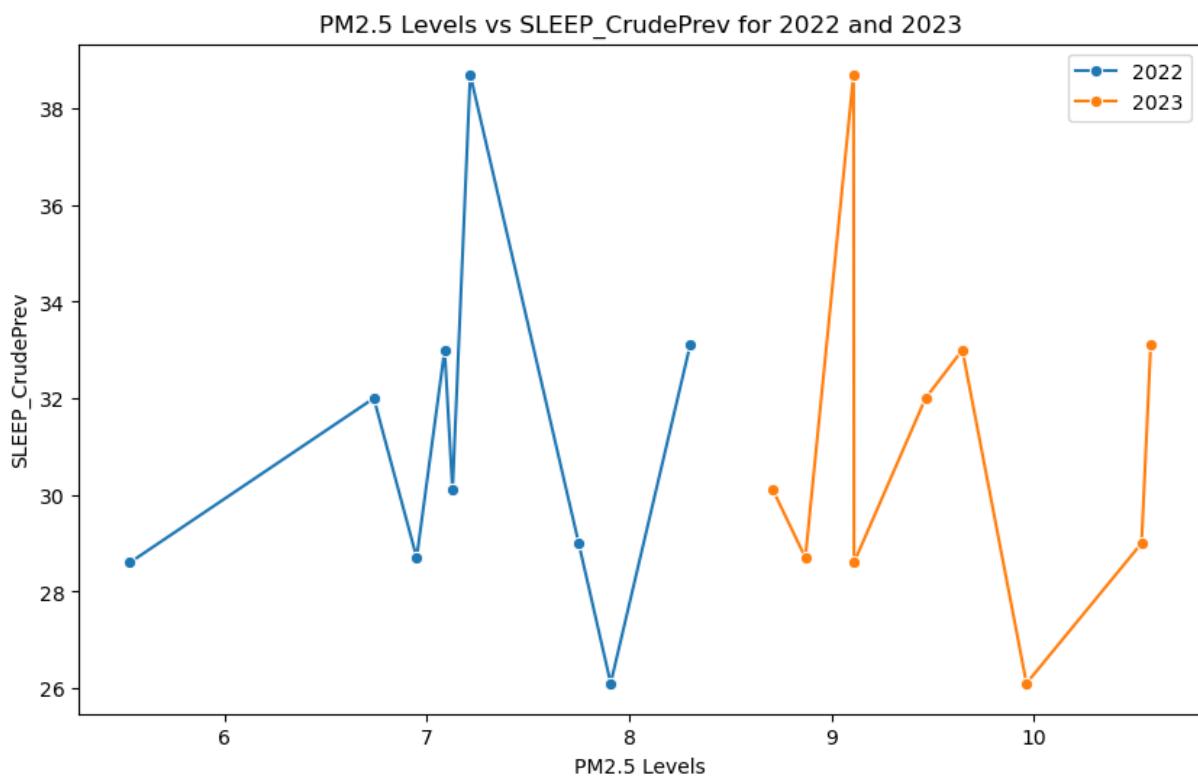
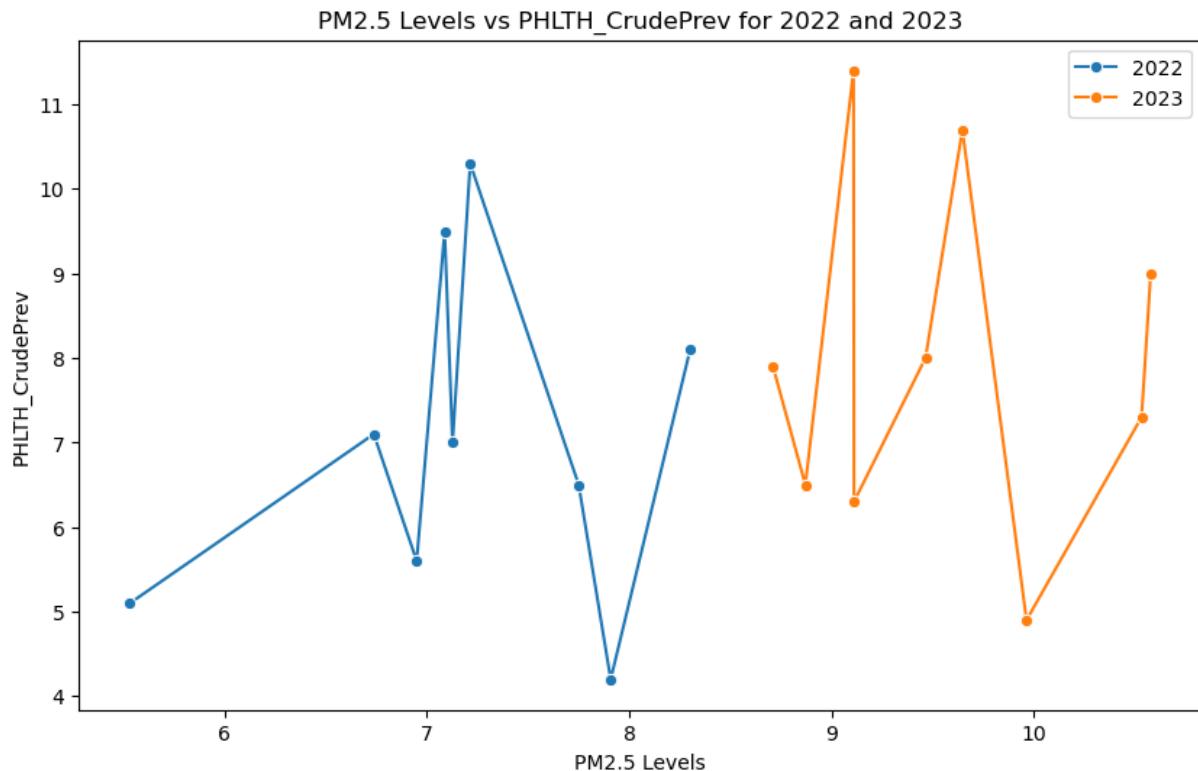


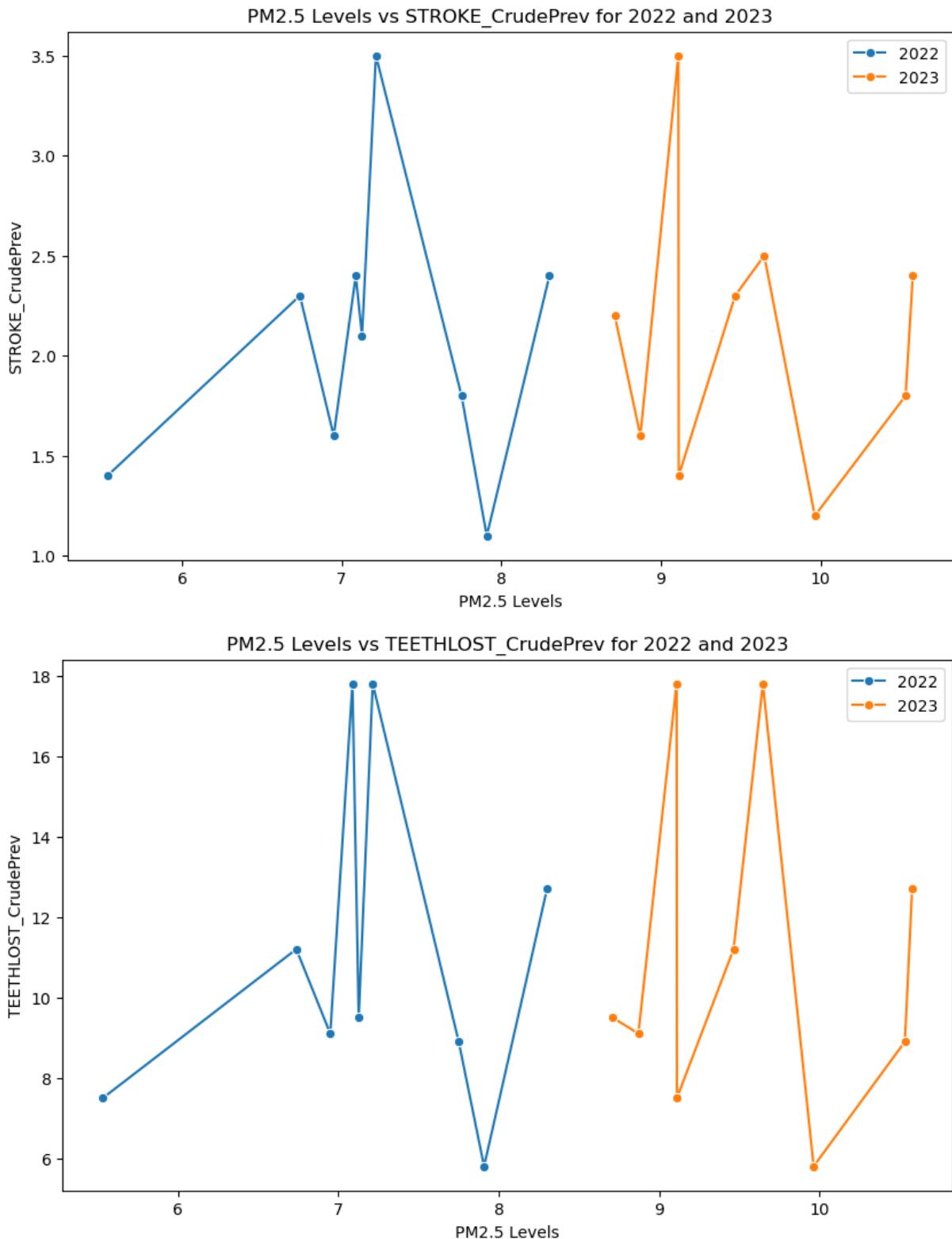












```
In [65]: health_data_2022['geometry'] = health_data_2022['Geolocation'].apply(wkt.loads)
health_data_2023['geometry'] = health_data_2023['Geolocation'].apply(wkt.loads)

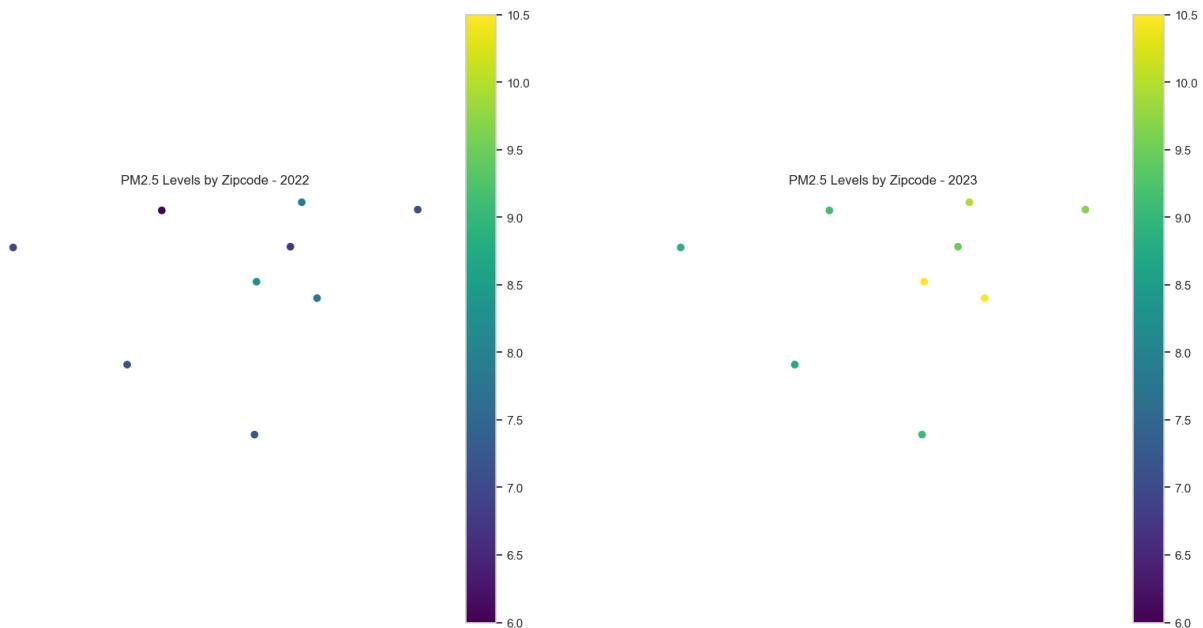
# Creating GeoDataFrame
geo_data_2022 = gpd.GeoDataFrame(health_data_2022, geometry='geometry')
geo_data_2023 = gpd.GeoDataFrame(health_data_2023, geometry='geometry')

fig, ax = plt.subplots(1, 2, figsize=(20, 10))
```

```
# Map for 2022
geo_data_2022.plot(column='PM2.5', cmap='viridis', legend=True, ax=ax[0], vrange=[6.0, 10.5])
ax[0].set_title('PM2.5 Levels by Zipcode - 2022')
ax[0].set_axis_off()

# Map for 2023
geo_data_2023.plot(column='PM2.5', cmap='viridis', legend=True, ax=ax[1], vrange=[6.0, 10.5])
ax[1].set_title('PM2.5 Levels by Zipcode - 2023')
ax[1].set_axis_off()

plt.show()
```



```
In [13]: cohort_data = pd.merge(health_data_2022, health_data_2023, on='Zipcode', suffixes=('_2022', '_2023'))

# Preparing data for heatmap
cohort_data_comparison = cohort_data.copy()
for col in health_data_2022.columns:
    if 'Crude' in col or col == 'PM2.5':
        cohort_data_comparison[f'{col}_Change'] = cohort_data_comparison[f'{col}_2023'] - cohort_data_comparison[f'{col}_2022']

# Function to plot bar graphs for a selected health indicator
def plot_bar_graph_for_crude_rate(selected_rate):
    # Extracting relevant data
    data = cohort_data_comparison[['Zipcode', f'{selected_rate}_Change', 'PM2.5_Change']]

    # Plotting
    fig, ax = plt.subplots(figsize=(15, 8))
    data.set_index('Zipcode').plot(kind='bar', ax=ax)
    ax.set_title(f'Change in {selected_rate} and PM2.5 Levels by Zipcode (2023 vs 2022)')
    ax.set_xlabel('Zipcode')
    ax.set_ylabel('Change')
    plt.xticks(rotation=45)
    plt.show()

# Dropdown menu for selecting the crude rate
crude_rates = [col.replace('_2022', '').replace('_Change', '') for col in cohort_data_comparison.columns if 'Crude' in col]
```

```
dropdown = Dropdown(options=crude_rates, description='Select Crude Rate:')

# Interactive widget
interact(plot_bar_graph_for_crude_rate, selected_rate=dropdown)

interactive(children=(Dropdown(description='Select Crude Rate:', options=('ACCESS2_CrudePrev', 'ARTHRITIS_Crud...')))

Out[13]: <function __main__.plot_bar_graph_for_crude_rate(selected_rate)>
```

Boston - Analysis

```
In [14]: boston_df = pd.read_csv('../Merged_Data/Boston_Health_20_23.csv')

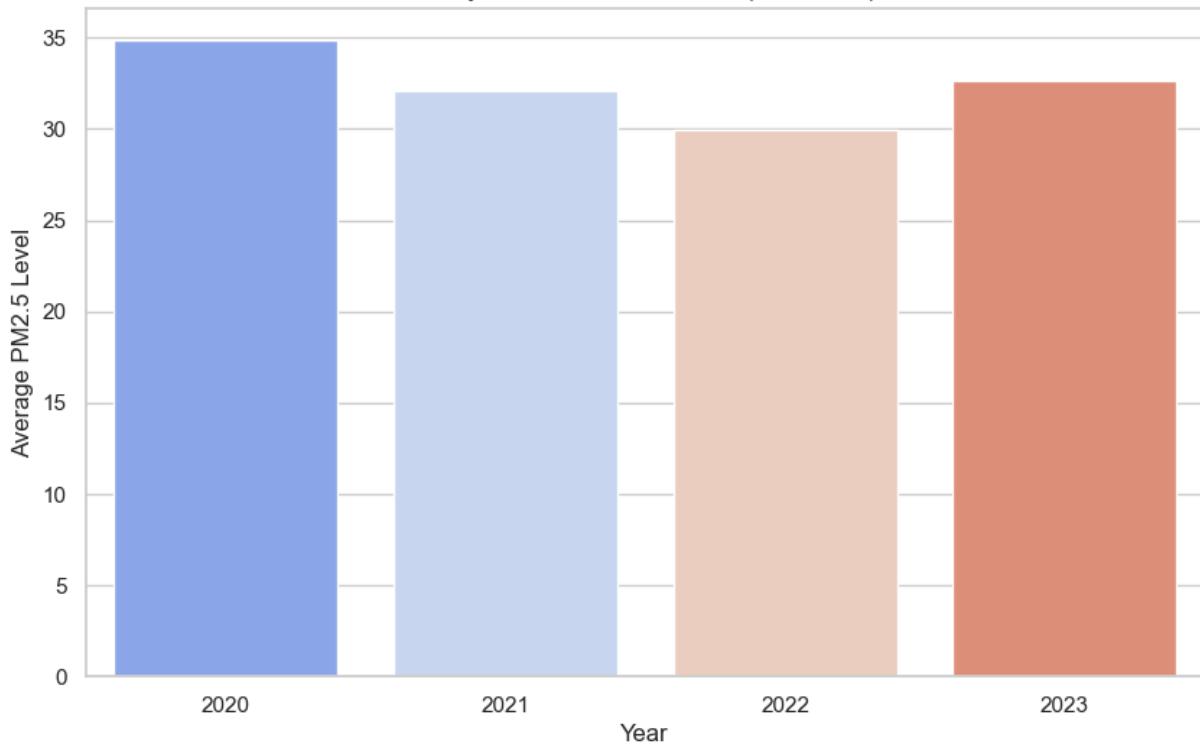
In [15]: sns.set(style="whitegrid")

# Plotting PM2.5 levels for each year
plt.figure(figsize=(10, 6))
sns.barplot(x="Year", y="pm25", data=boston_df, palette="coolwarm")
plt.title("Yearly PM2.5 Levels in Boston (2020–2023)")
plt.xlabel("Year")
plt.ylabel("Average PM2.5 Level")
plt.show()

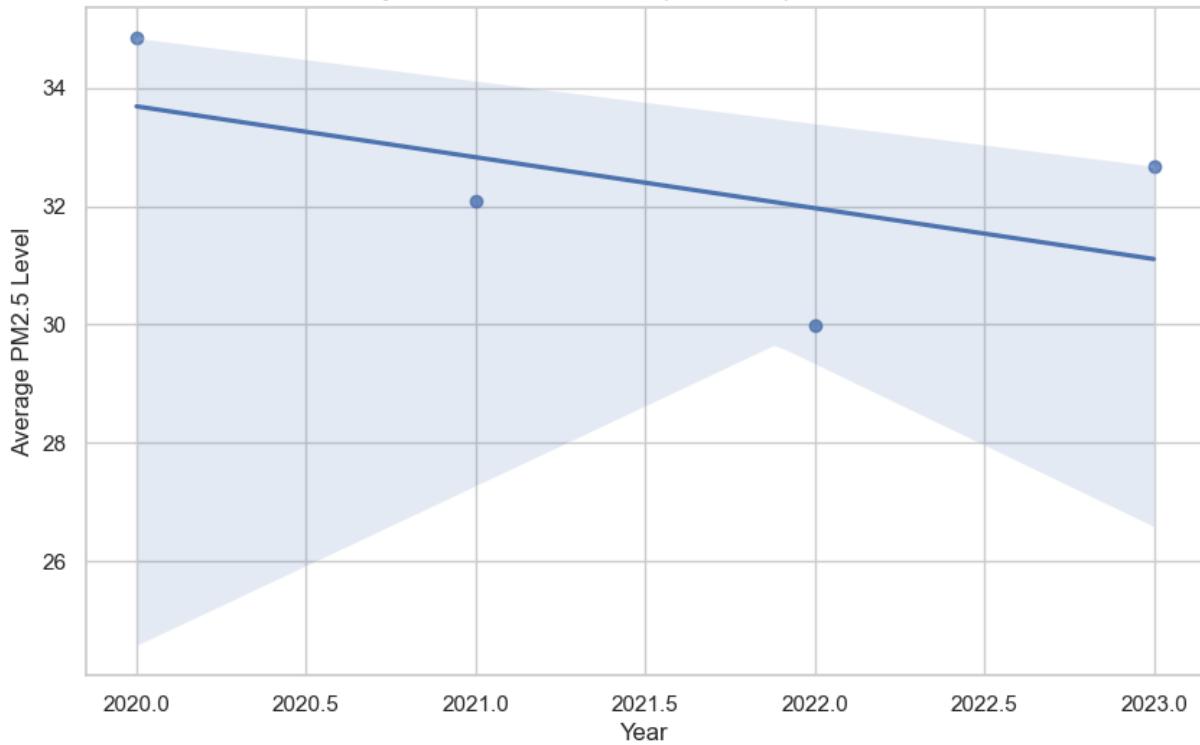
plt.figure(figsize=(10, 6))
sns.regplot(x="Year", y="pm25", data=boston_df, color='b', marker='o')

# Enhancing the plot
plt.title("Yearly PM2.5 Levels in Boston (2020–2023) with Trend Line")
plt.xlabel("Year")
plt.ylabel("Average PM2.5 Level")
plt.show()
```

Yearly PM2.5 Levels in Boston (2020-2023)



Yearly PM2.5 Levels in Boston (2020-2023) with Trend Line

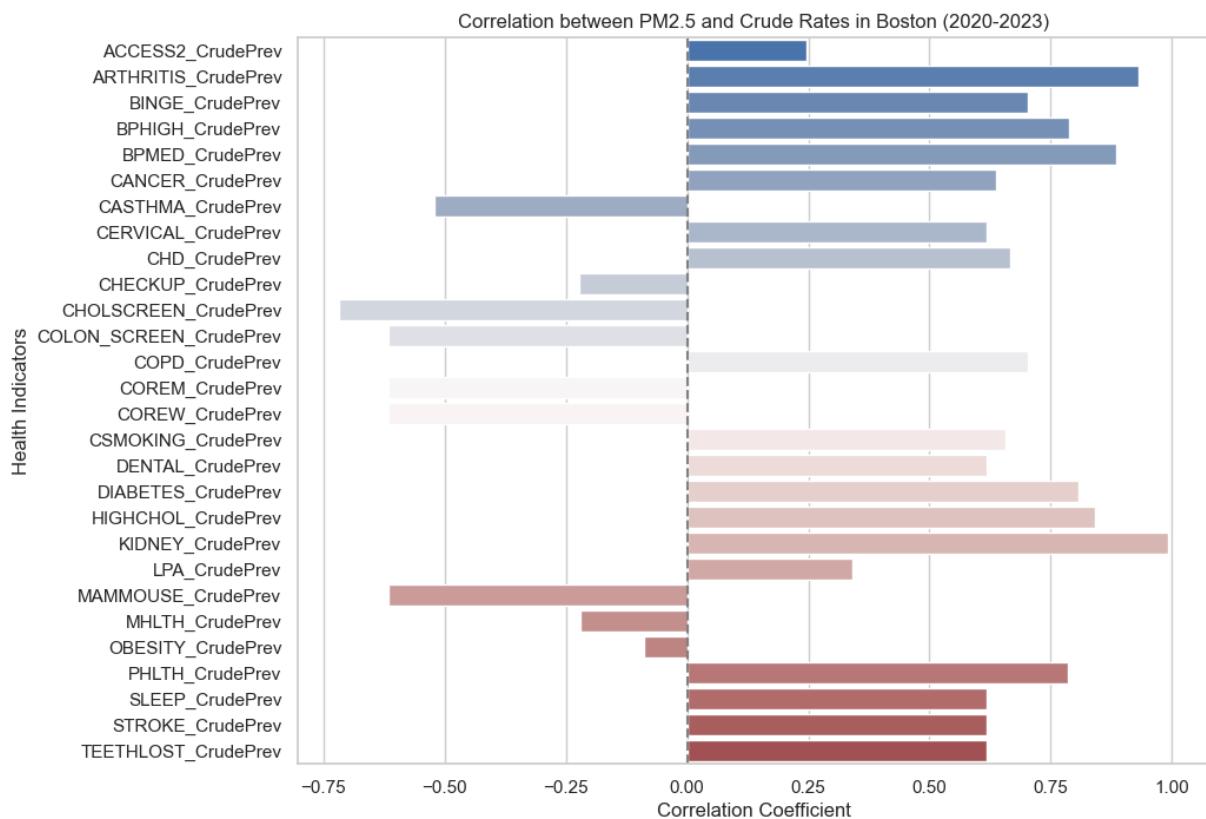


```
In [16]: crude_columns = [col for col in boston_df.columns if 'Crude' in col]
crude_columns.append('pm25')

# Calculating the correlation matrix for the specified columns in boston_df
correlation_matrix = boston_df[crude_columns].corr()

# Extracting correlations with pm25, excluding the self-correlation of pm25
correlations_pm25 = correlation_matrix['pm25'].drop('pm25')
```

```
# Creating a bar plot to visualize the correlations
plt.figure(figsize=(10, 8))
sns.barplot(x=correlations_pm25.values, y=correlations_pm25.index, palette="vlag")
plt.title('Correlation between PM2.5 and Crude Rates in Boston (2020-2023)')
plt.xlabel('Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```



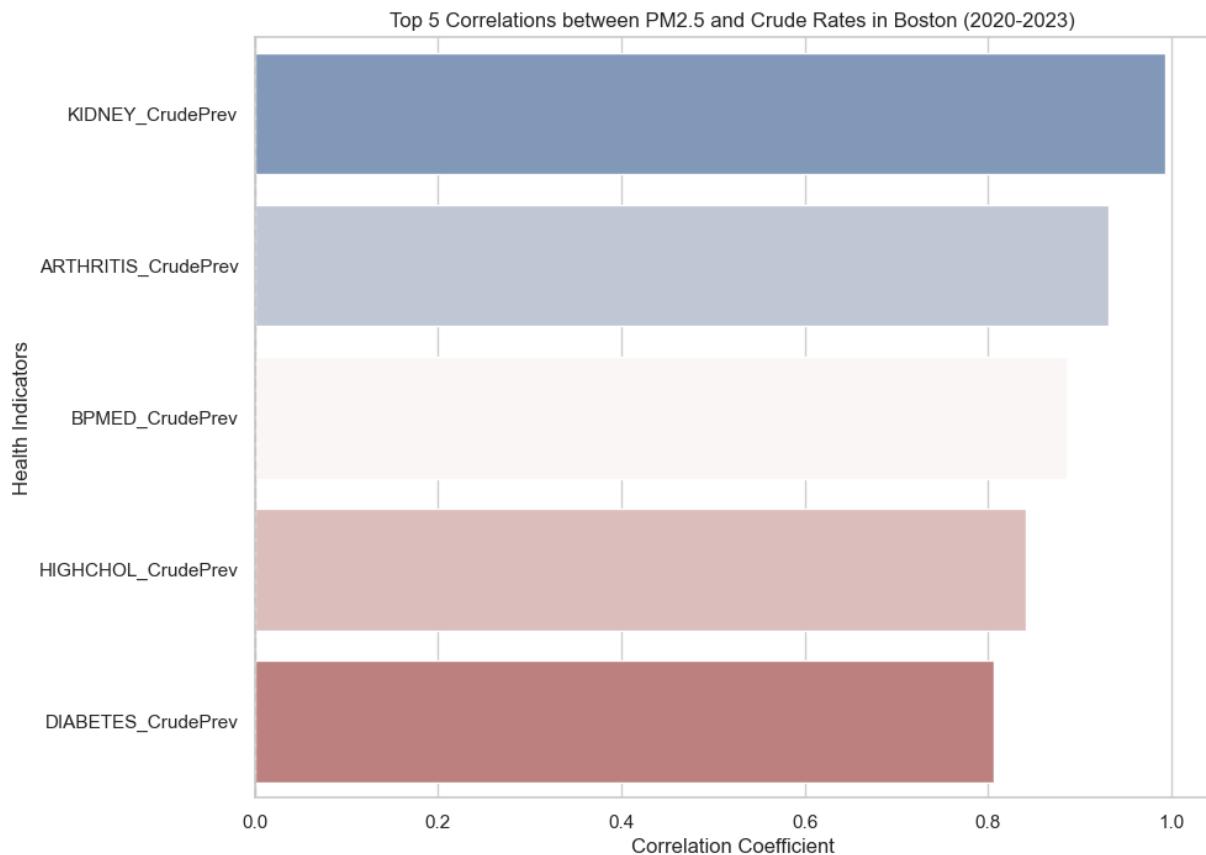
```
In [68]: import matplotlib.pyplot as plt
import seaborn as sns

# Assuming you've already calculated 'correlations_pm25' as in your provided code cell above.

# Sort correlations in descending order and get the top 5
top_correlations = correlations_pm25.abs().sort_values(ascending=False).head(5)

# Subset correlation matrix to the top 5 correlated columns with PM2.5
top_correlation_matrix = correlation_matrix.loc[top_correlations.index, 'pm25']

# Create a bar plot for the top 5 correlations
plt.figure(figsize=(10, 8))
sns.barplot(x=top_correlations.values, y=top_correlations.index, palette="vlag")
plt.title('Top 5 Correlations between PM2.5 and Crude Rates in Boston (2020-2023)')
plt.xlabel('Correlation Coefficient')
plt.ylabel('Health Indicators')
plt.axvline(x=0, color='gray', linestyle='--')
plt.show()
```



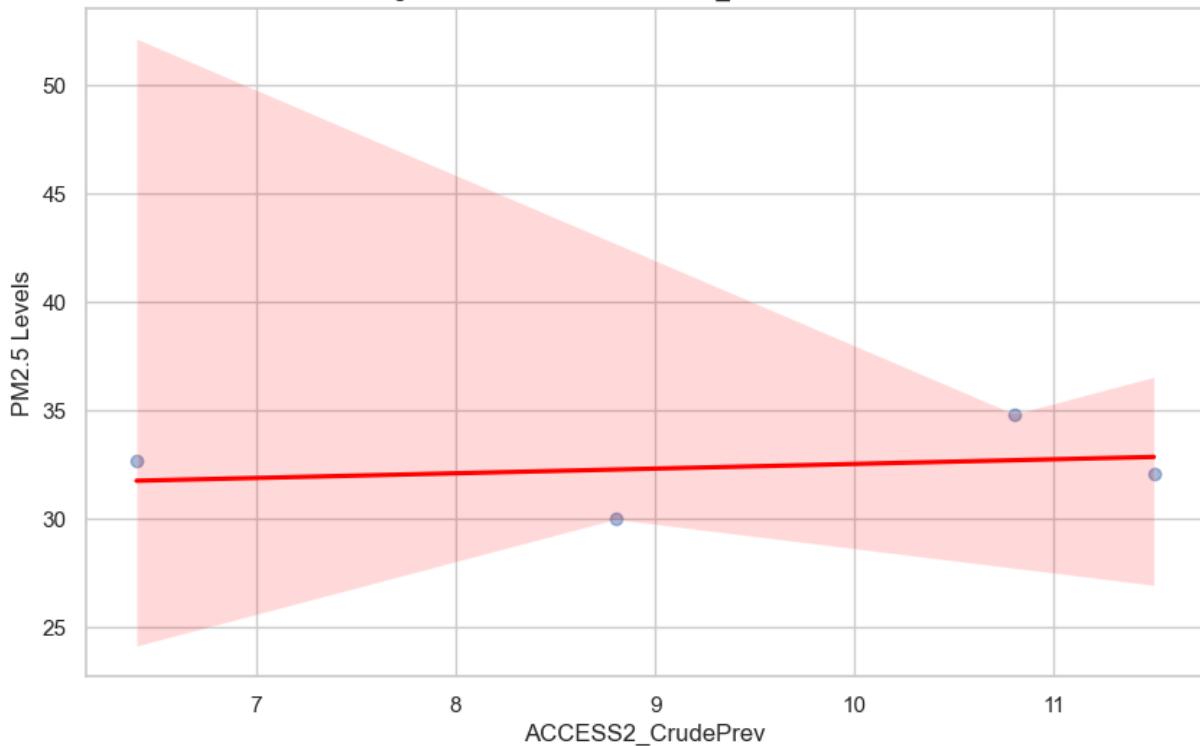
```
In [17]: def plot_linear_regression(data, crude_columns, year):
    for column in crude_columns:
        if column != 'pm25':
            # Performing linear regression
            slope, intercept, r_value, p_value, std_err = stats.linregress(data['pm25'], data[column])

            # Plotting
            plt.figure(figsize=(10, 6))
            sns.regplot(x=column, y='pm25', data=data, scatter_kws={'alpha': 0.5})
            plt.title(f'Linear Regression: PM2.5 vs {column} for {year}')
            plt.xlabel(column)
            plt.ylabel('PM2.5 Levels')
            plt.show()

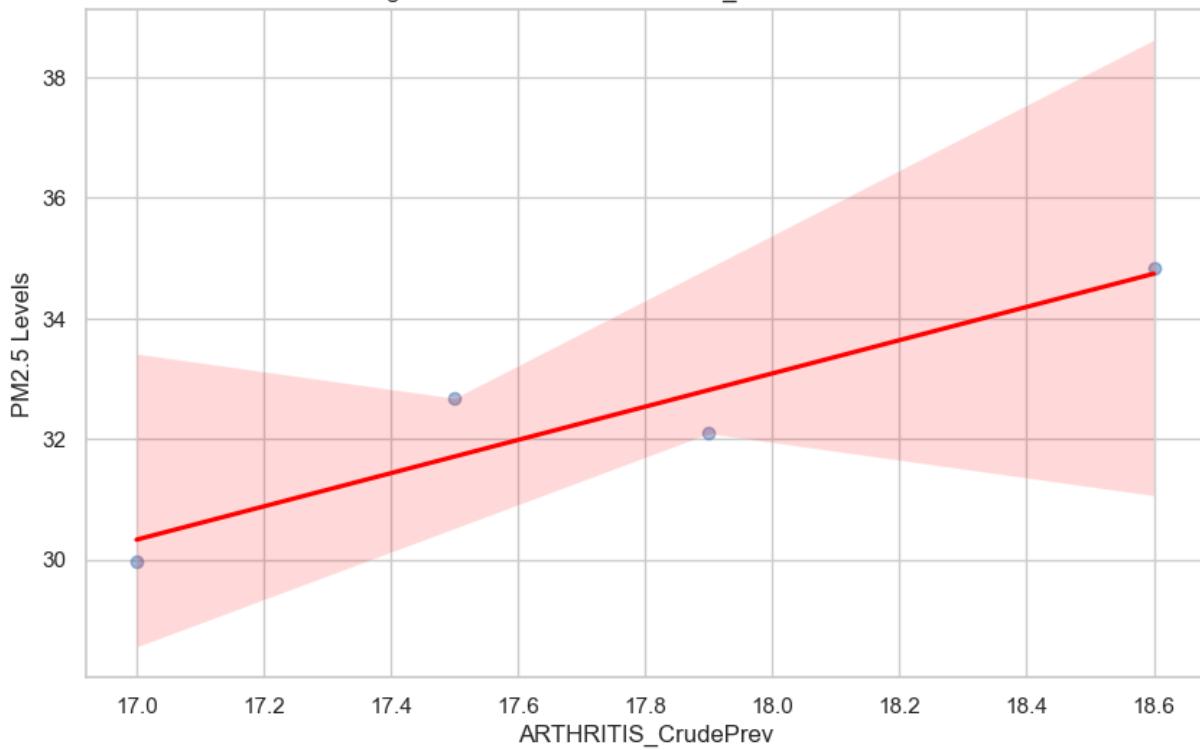
    # Extracting crude columns from the boston_df dataset
    crude_columns_boston = [col for col in boston_df.columns if 'Crude' in col]

    # Making linear regression plots for each crude rate against PM2.5 level for
    plot_linear_regression(boston_df, crude_columns_boston, '2020-2023')
```

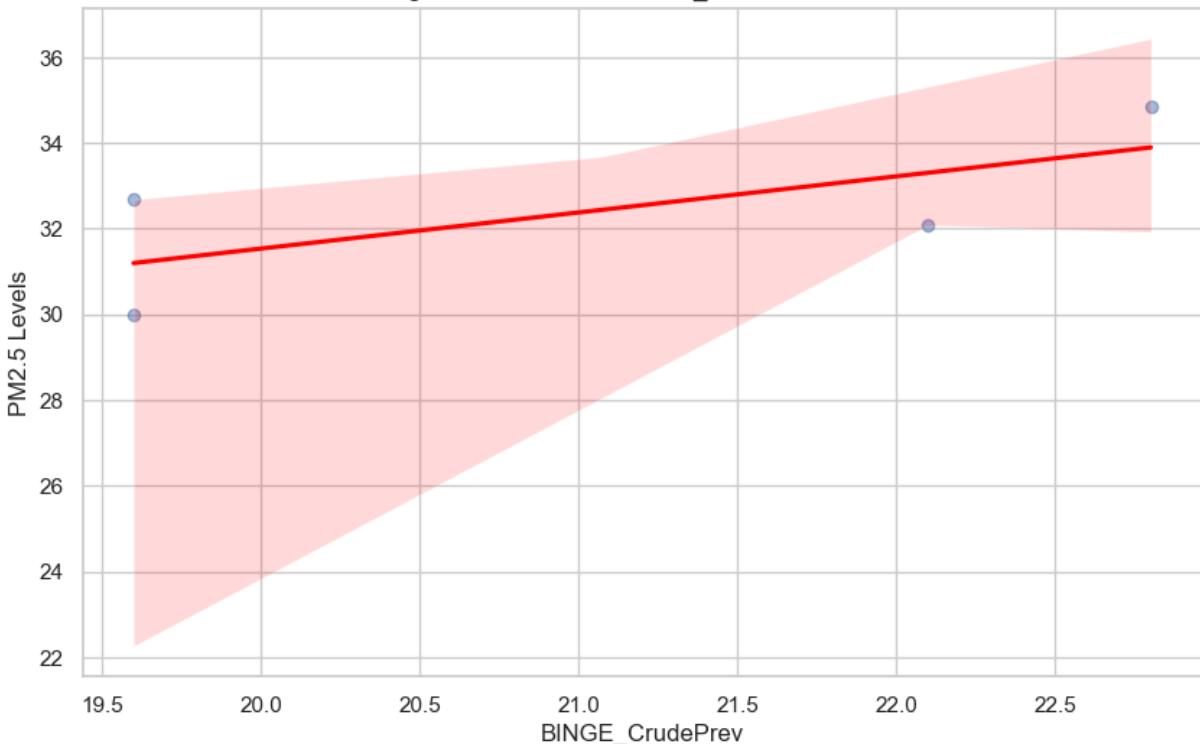
Linear Regression: PM2.5 vs ACCESS2_CrudePrev for 2020-2023



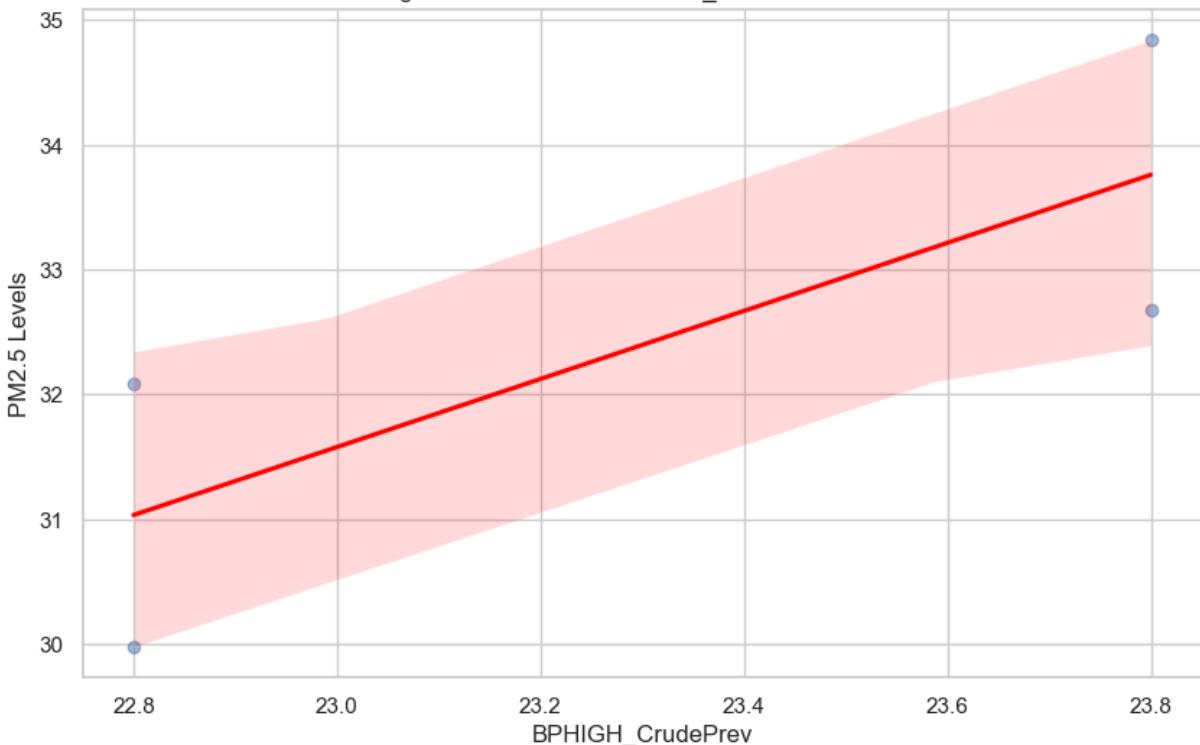
Linear Regression: PM2.5 vs ARTHRITIS_CrudePrev for 2020-2023



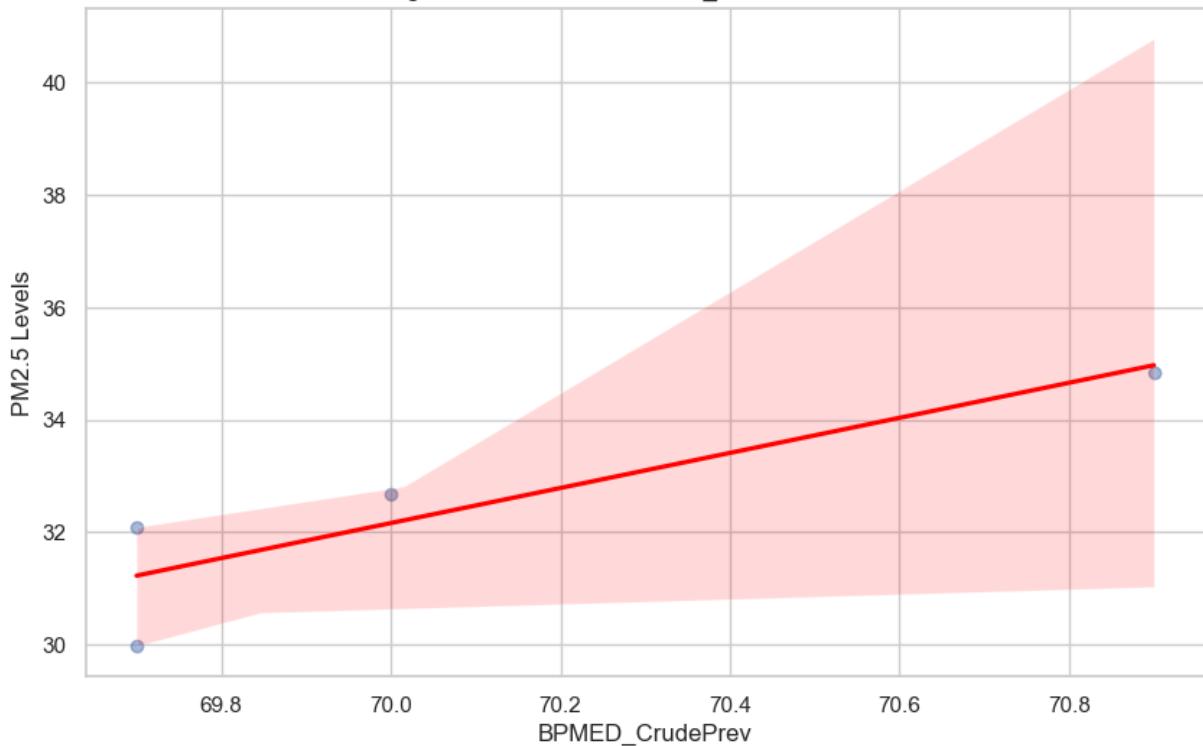
Linear Regression: PM2.5 vs BINGE_CrudePrev for 2020-2023



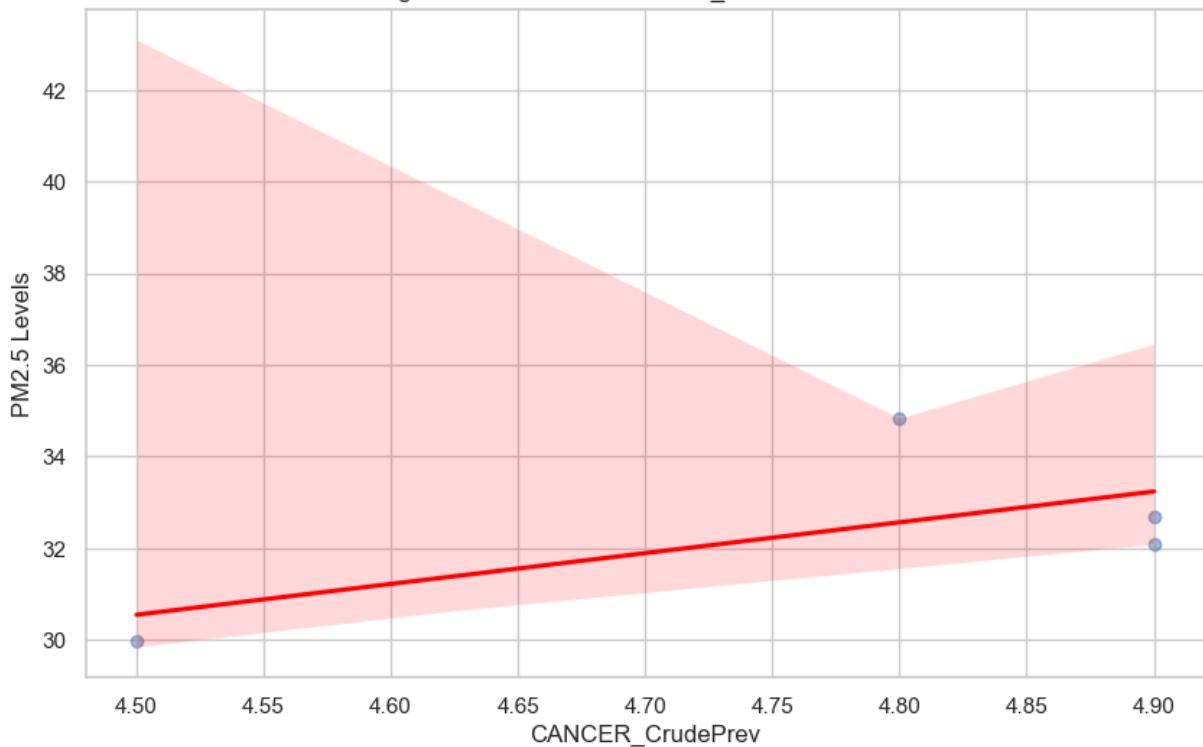
Linear Regression: PM2.5 vs BPHIGH_CrudePrev for 2020-2023



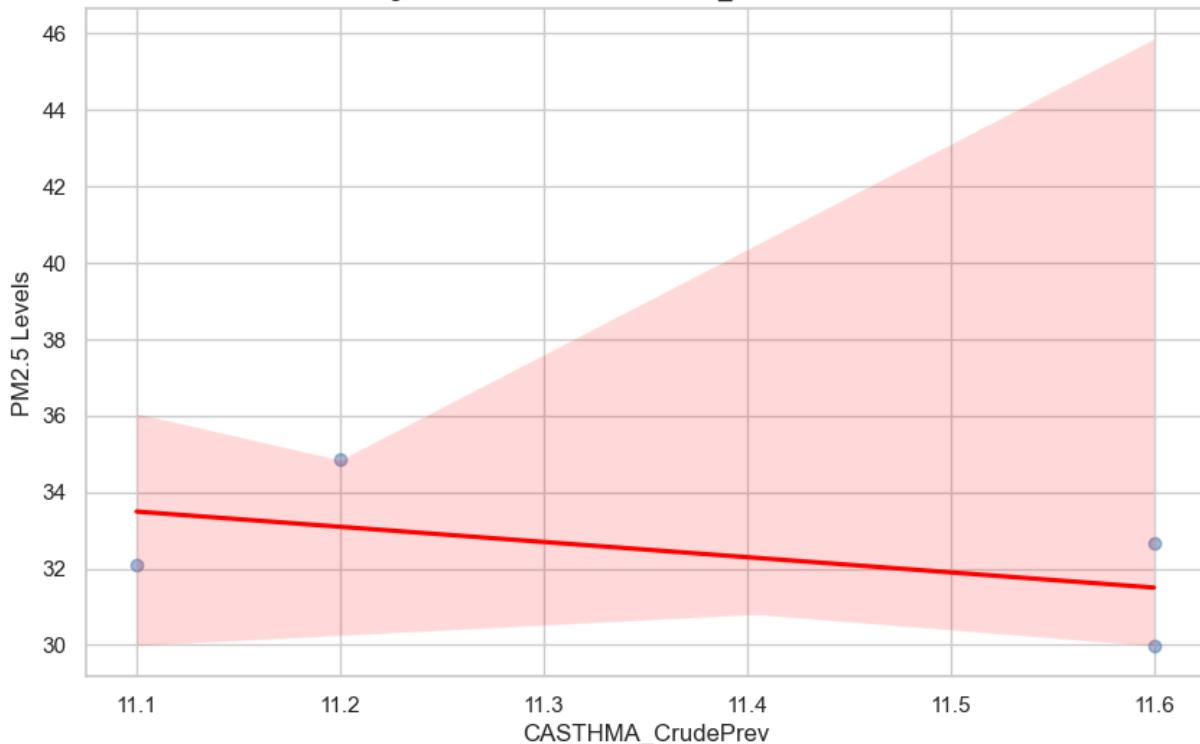
Linear Regression: PM2.5 vs BPMED_CrudePrev for 2020-2023



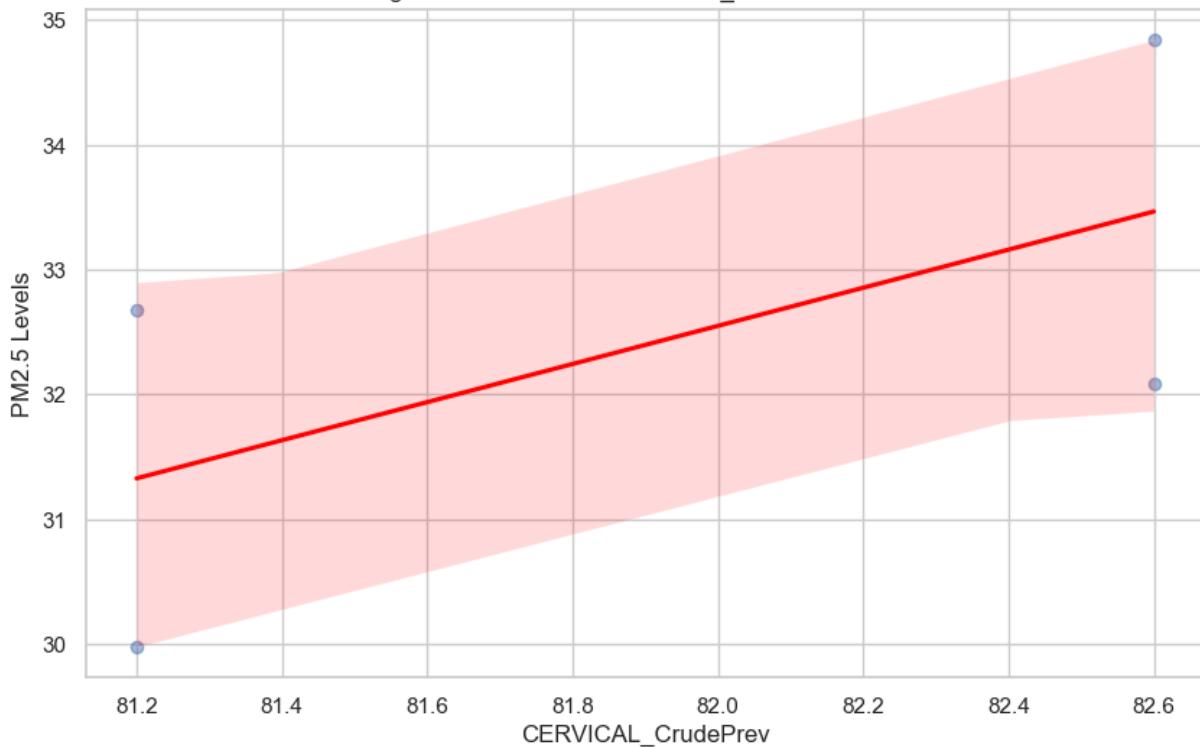
Linear Regression: PM2.5 vs CANCER_CrudePrev for 2020-2023



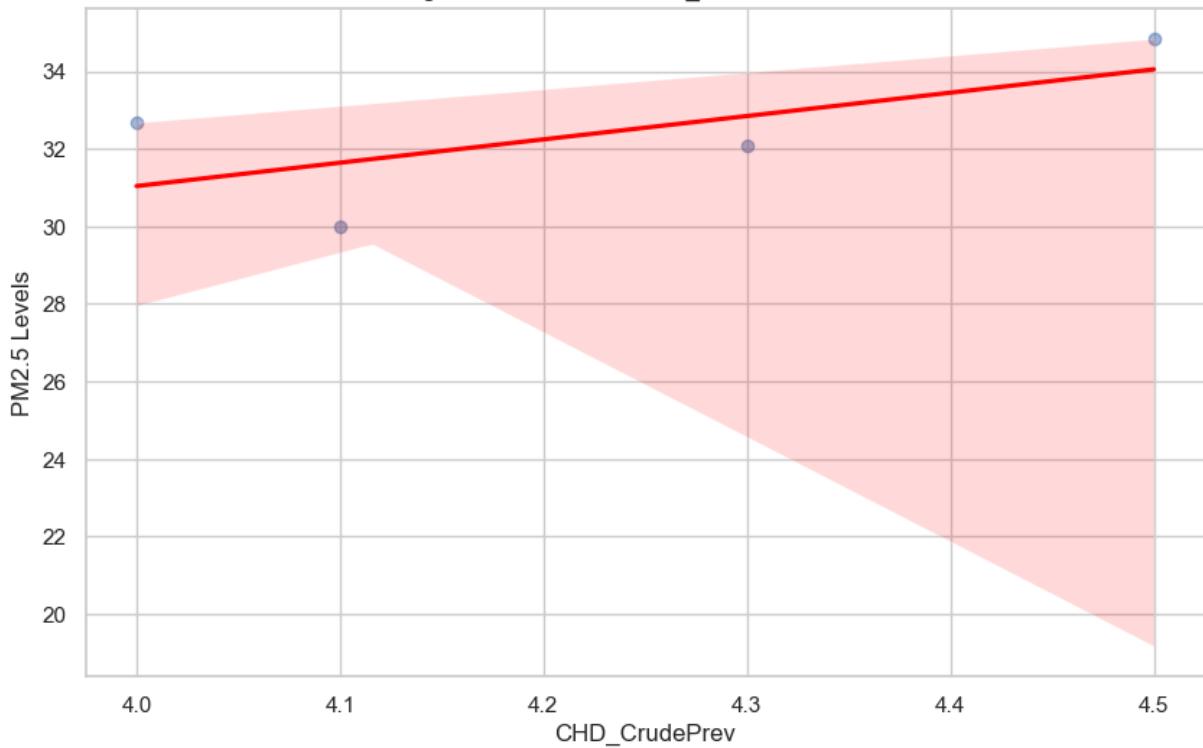
Linear Regression: PM2.5 vs CASTHMA_CrudePrev for 2020-2023



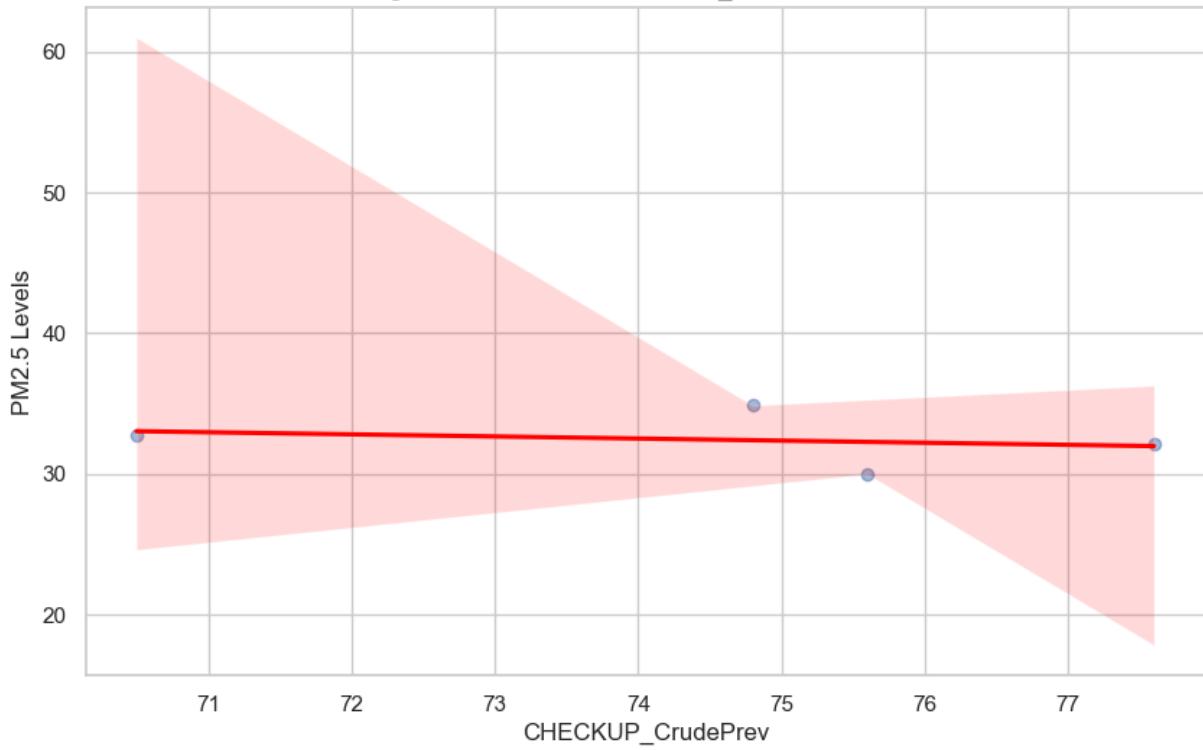
Linear Regression: PM2.5 vs CERVICAL_CrudePrev for 2020-2023



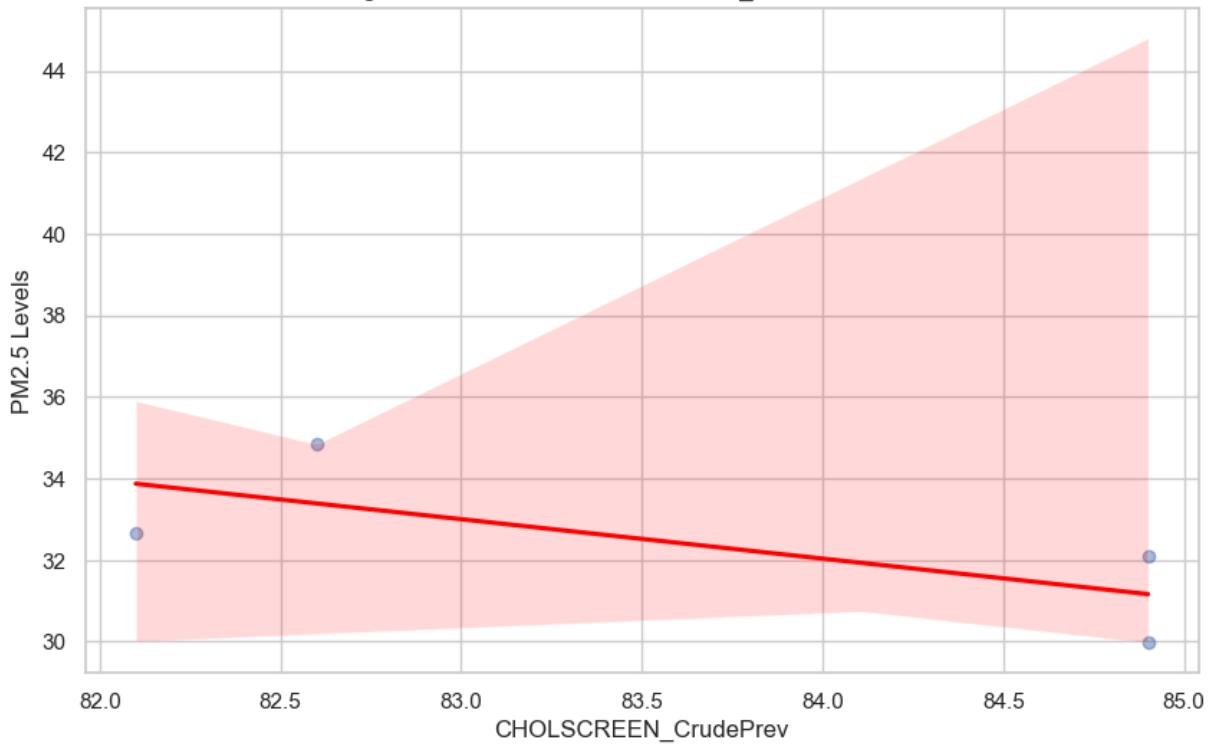
Linear Regression: PM2.5 vs CHD_CrudePrev for 2020-2023



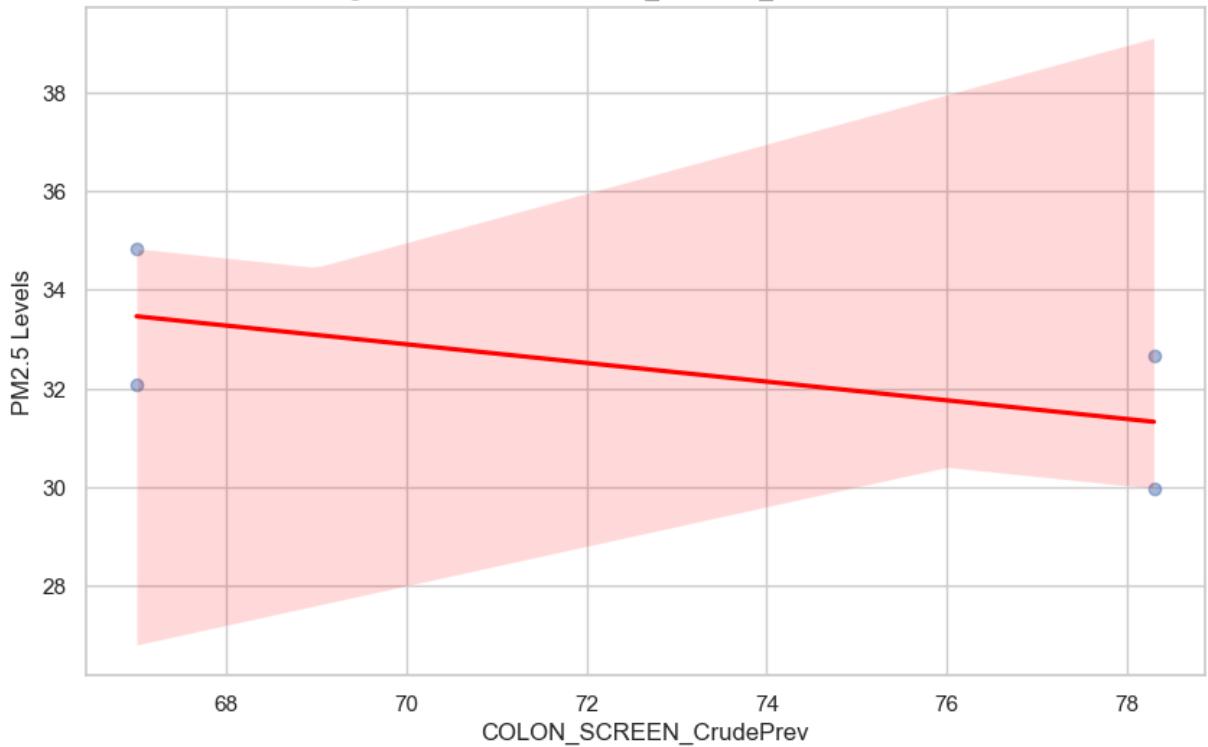
Linear Regression: PM2.5 vs CHECKUP_CrudePrev for 2020-2023

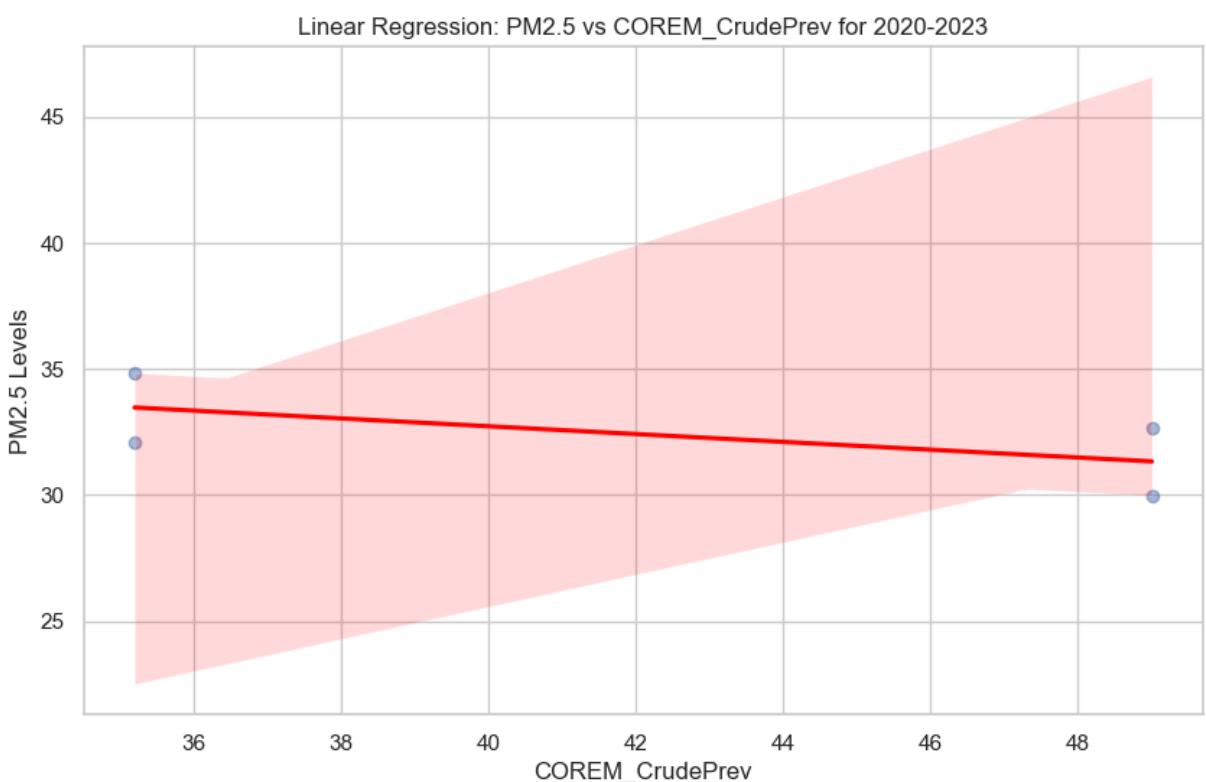


Linear Regression: PM2.5 vs CHOLSCREEN_CrudePrev for 2020-2023

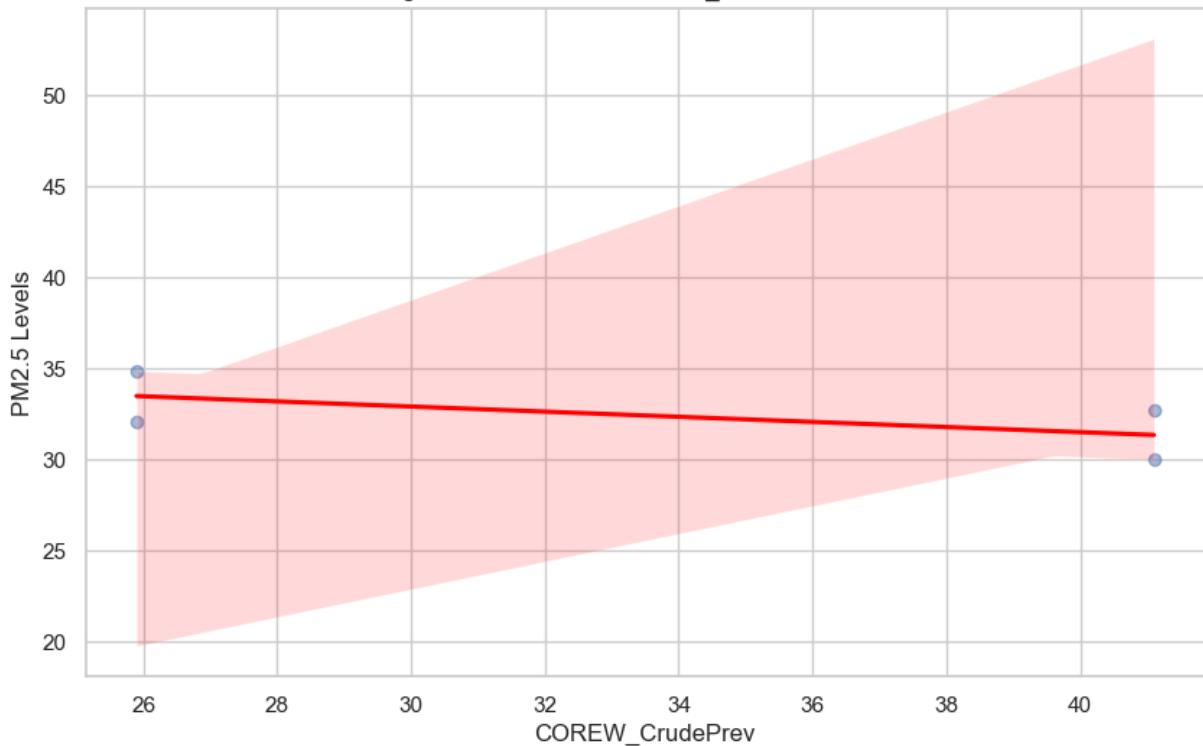


Linear Regression: PM2.5 vs COLON_SCREEN_CrudePrev for 2020-2023

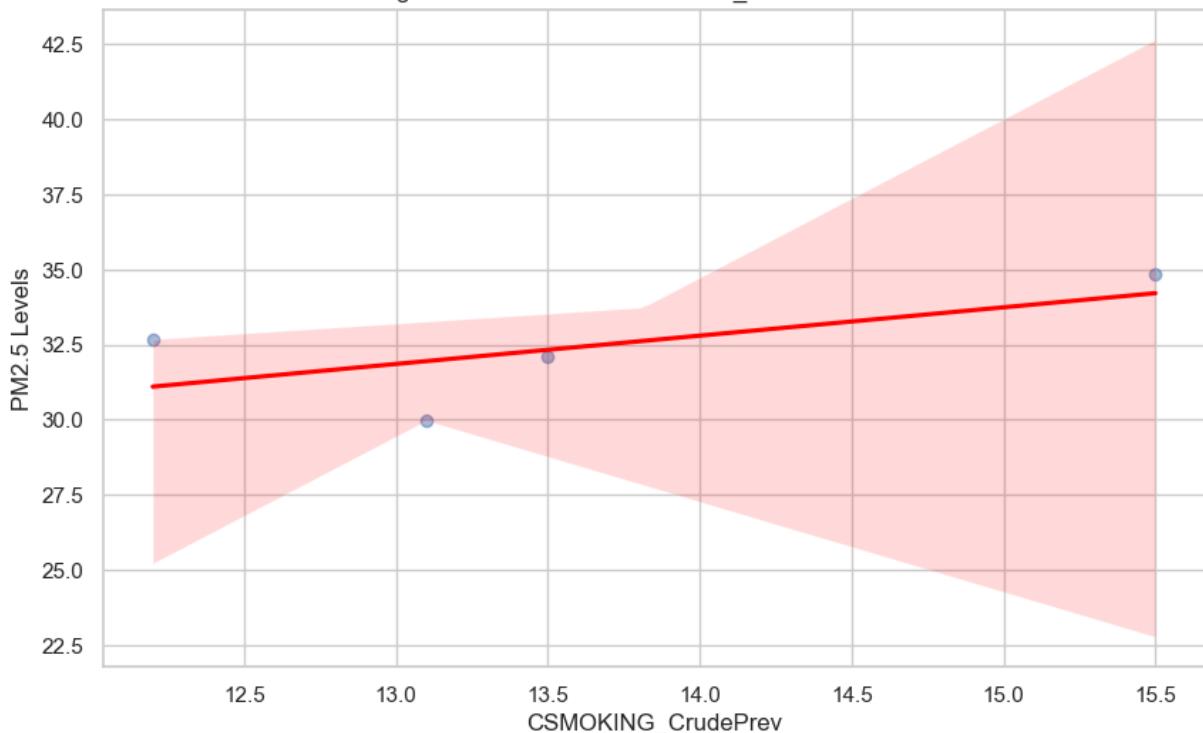




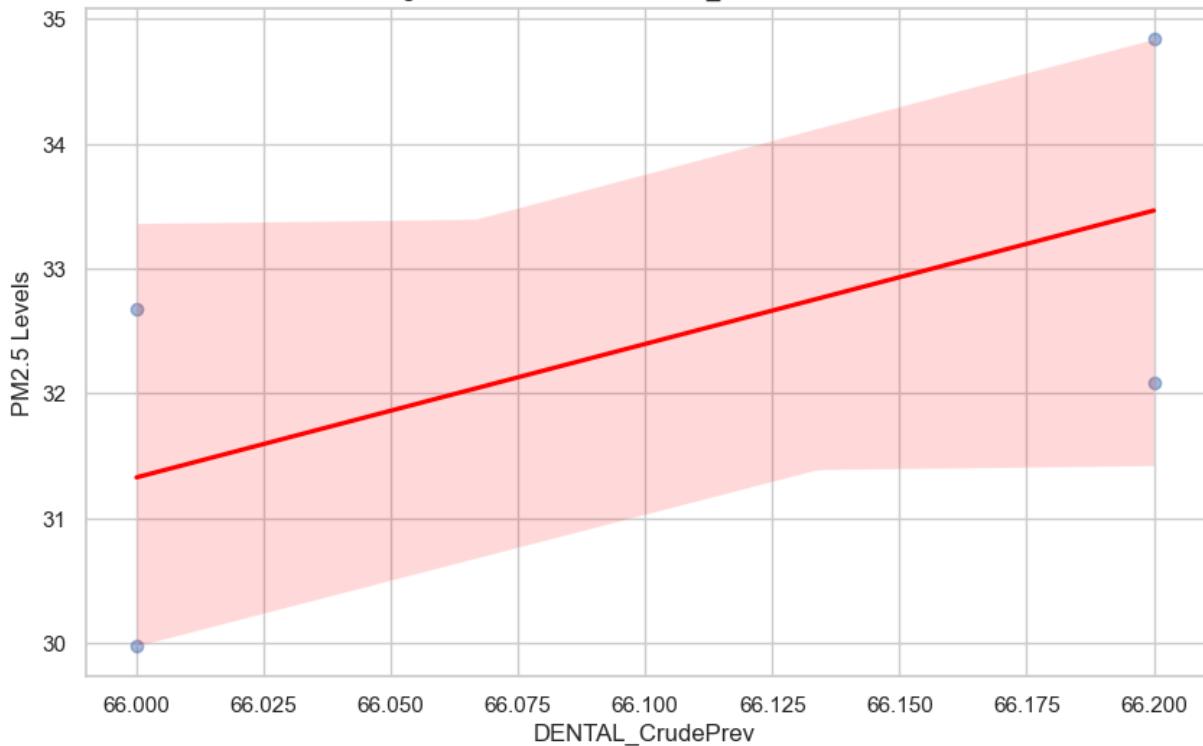
Linear Regression: PM2.5 vs COREW_CrudePrev for 2020-2023



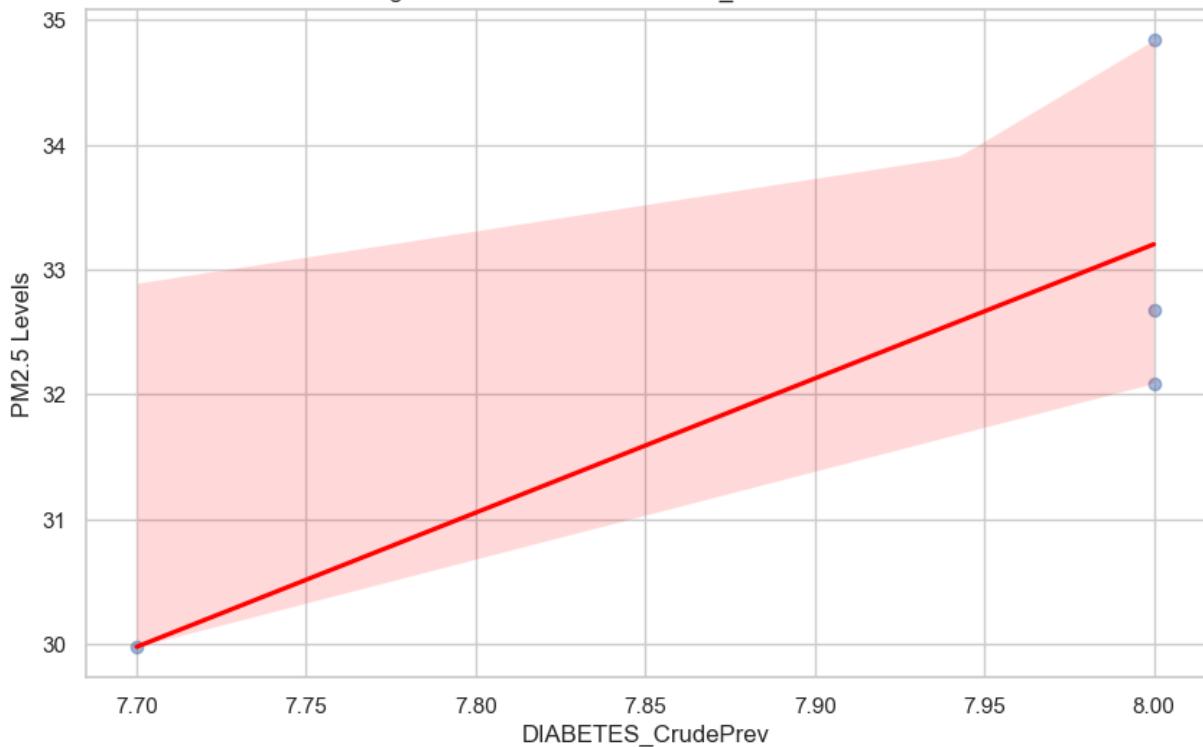
Linear Regression: PM2.5 vs CSMOKING_CrudePrev for 2020-2023

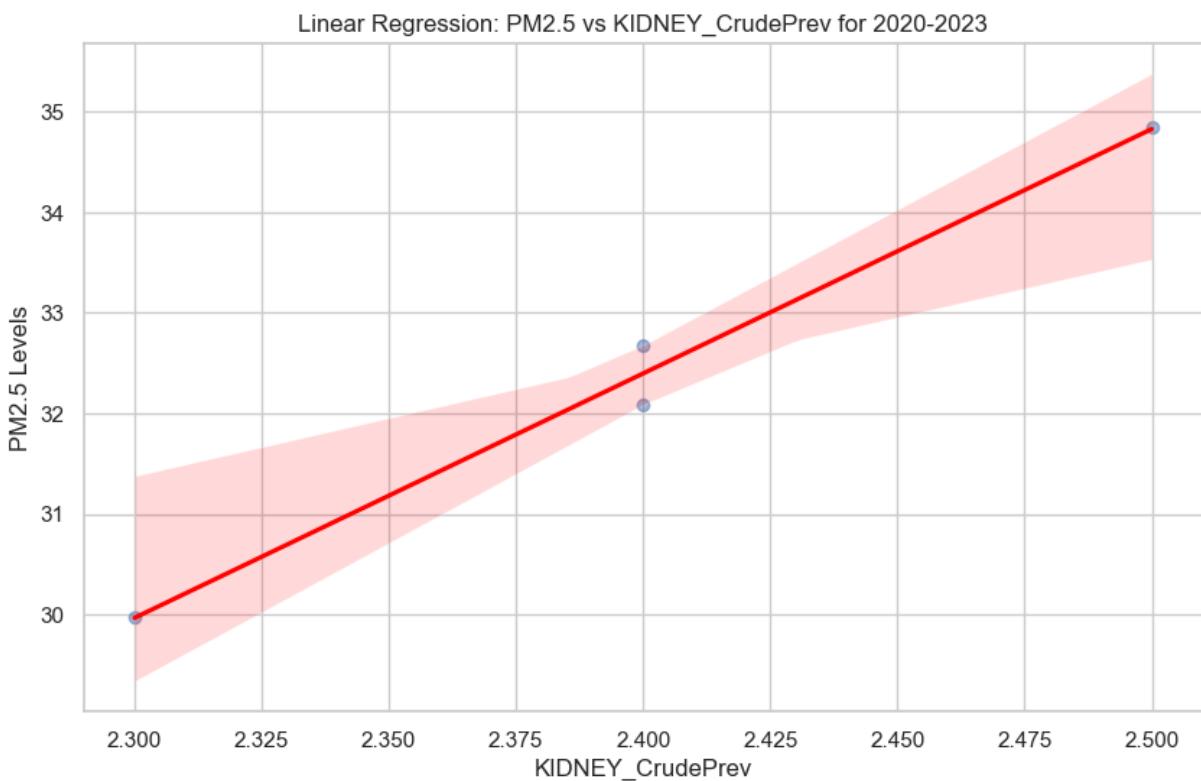
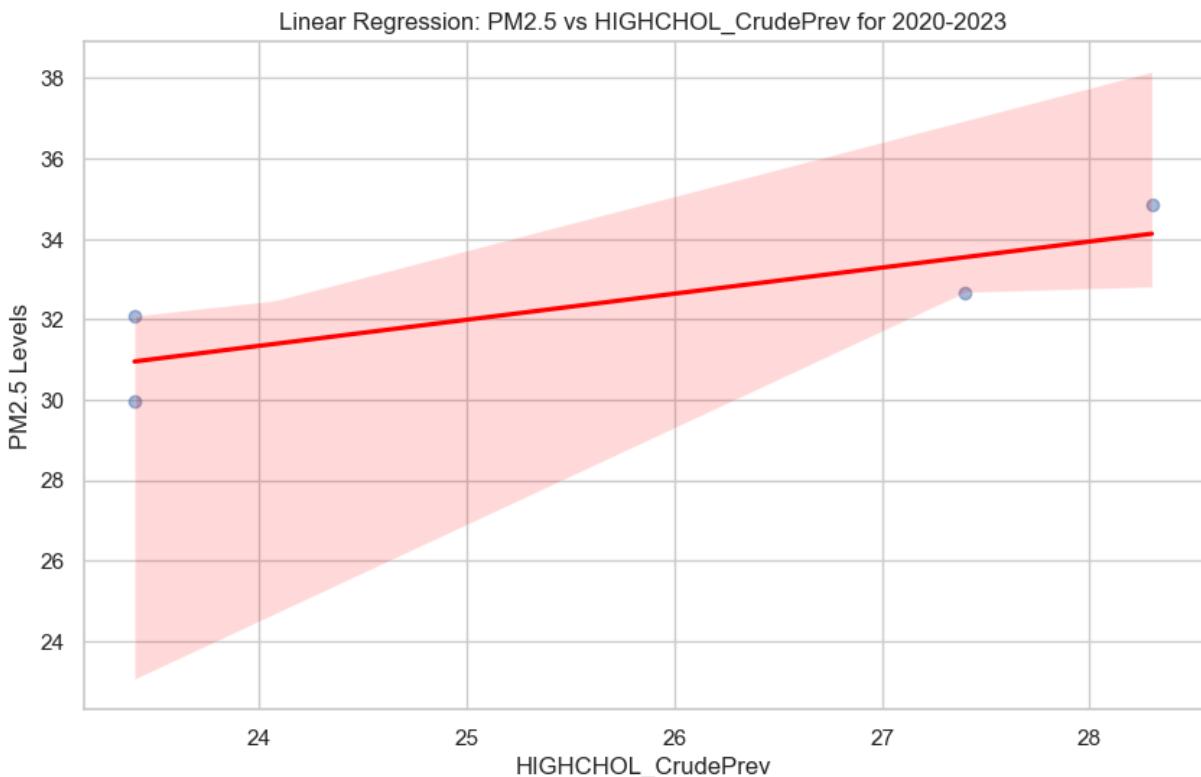


Linear Regression: PM2.5 vs DENTAL_CrudePrev for 2020-2023

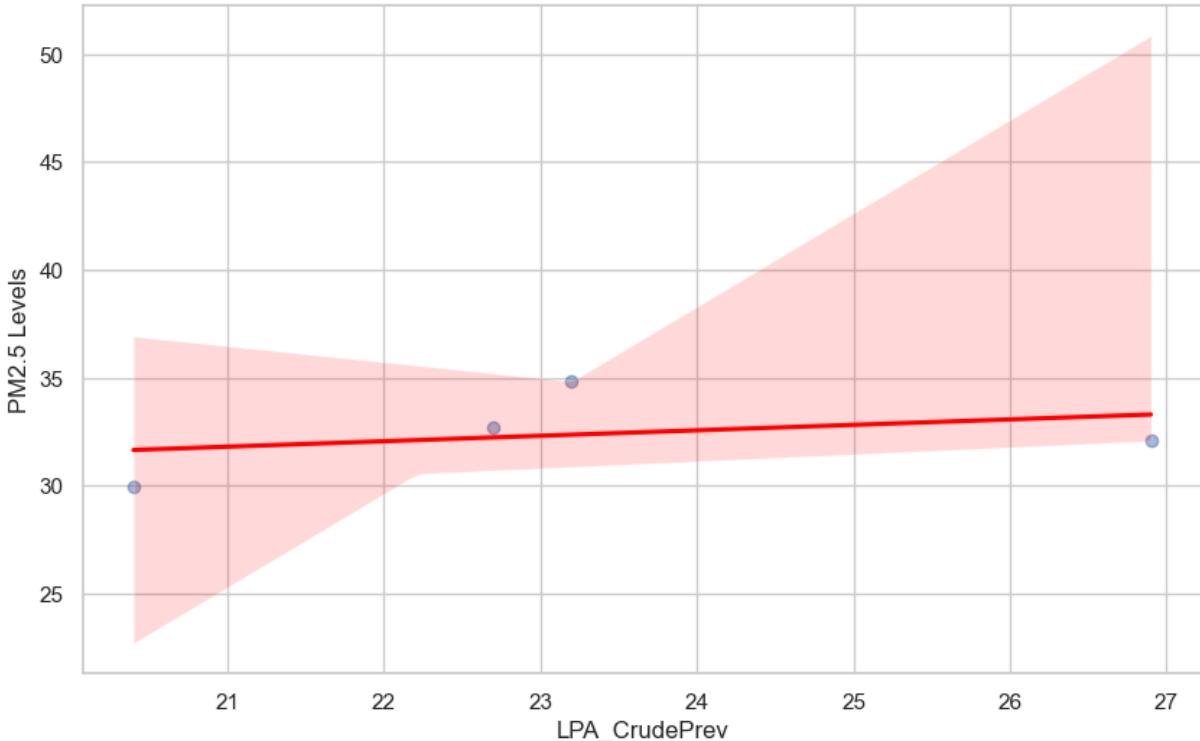


Linear Regression: PM2.5 vs DIABETES_CrudePrev for 2020-2023

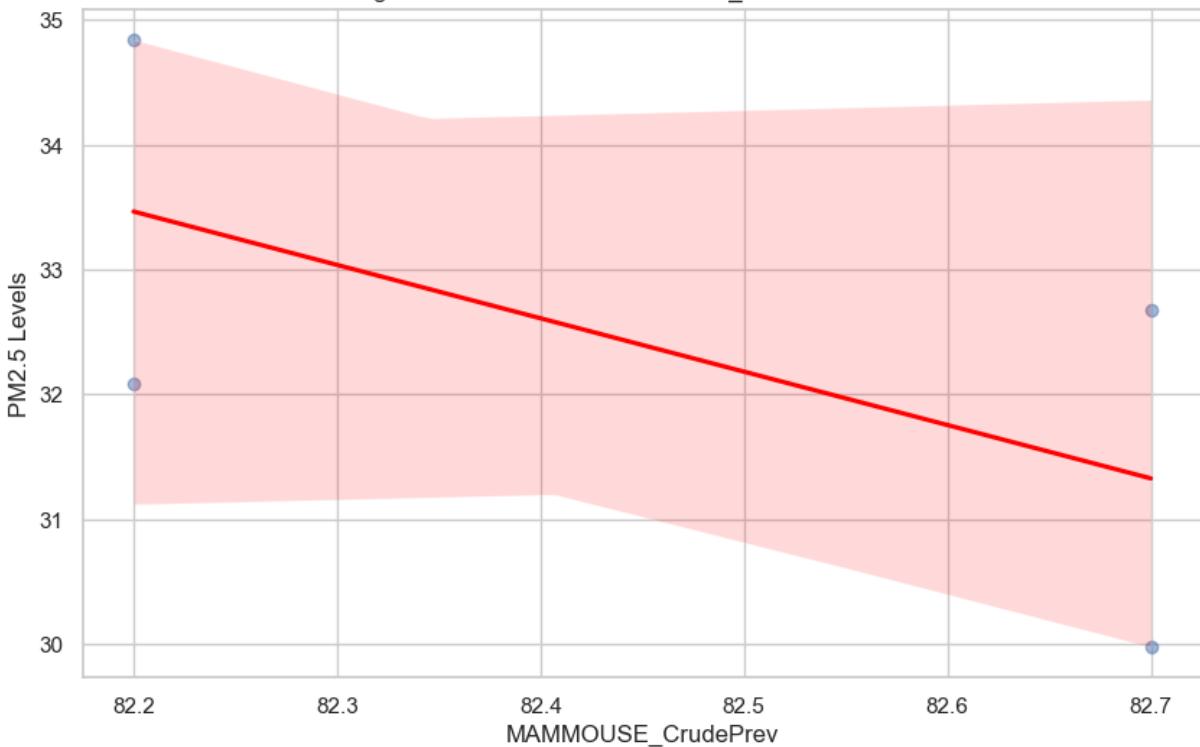




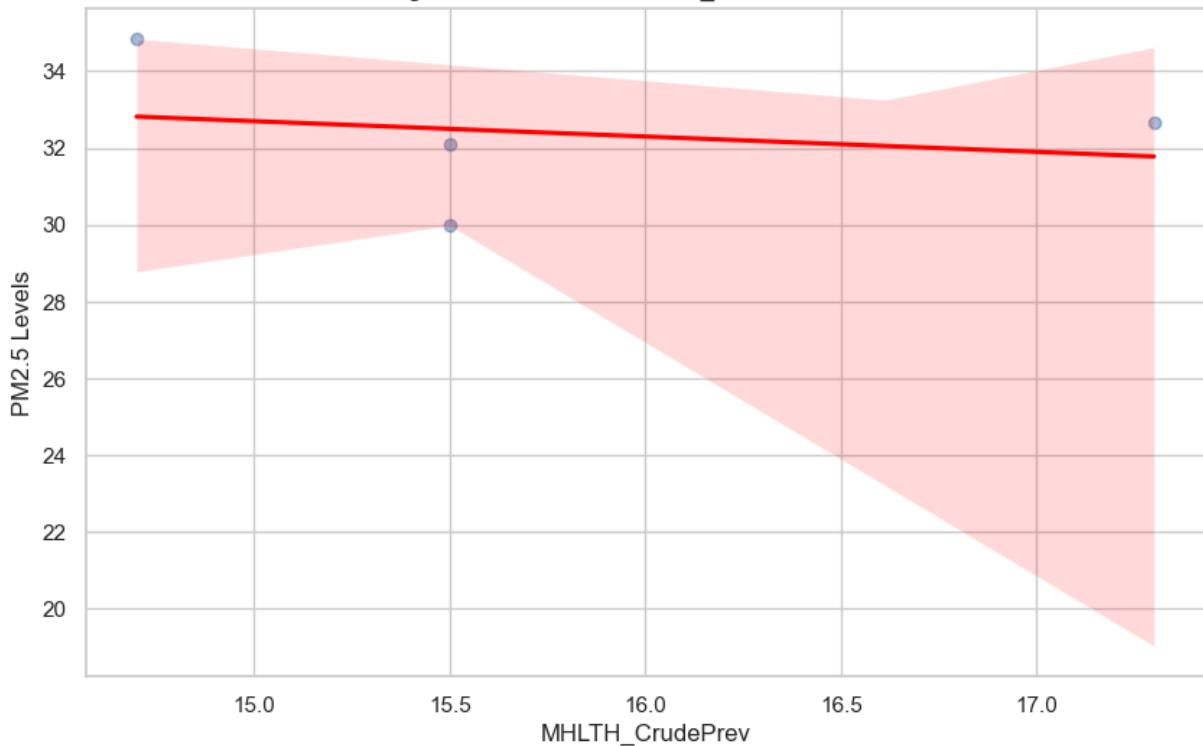
Linear Regression: PM2.5 vs LPA_CrudePrev for 2020-2023



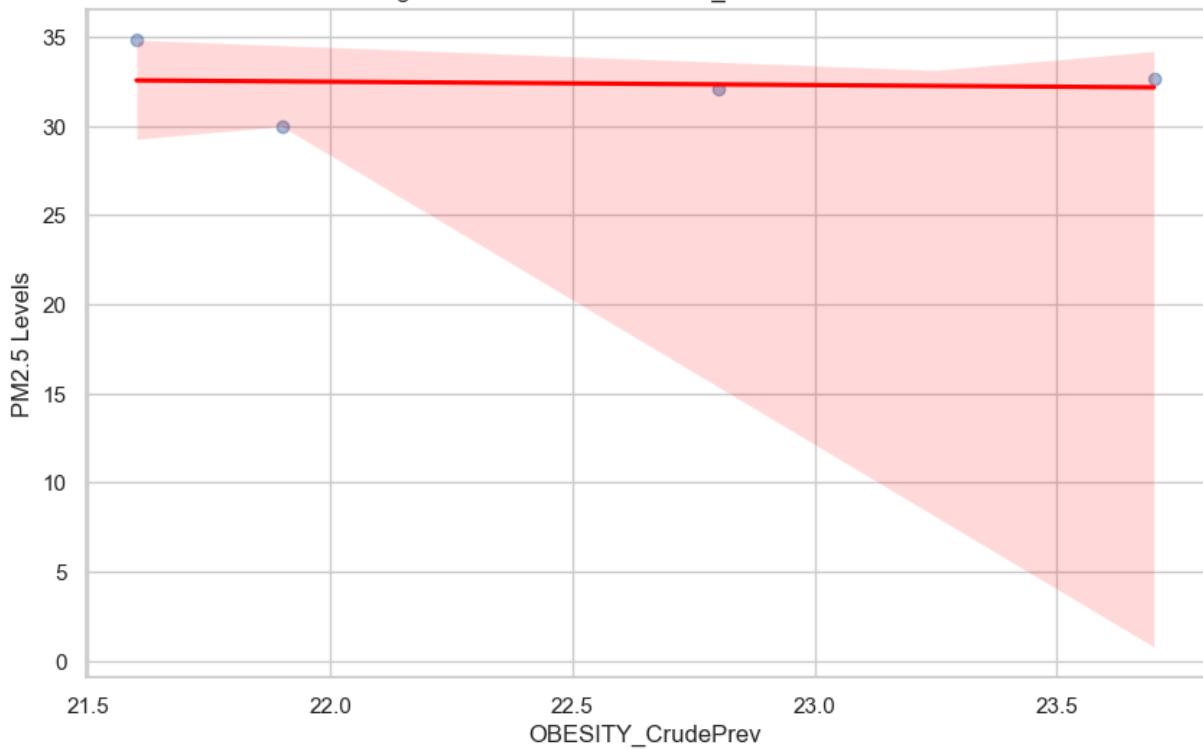
Linear Regression: PM2.5 vs MAMMOUSE_CrudePrev for 2020-2023

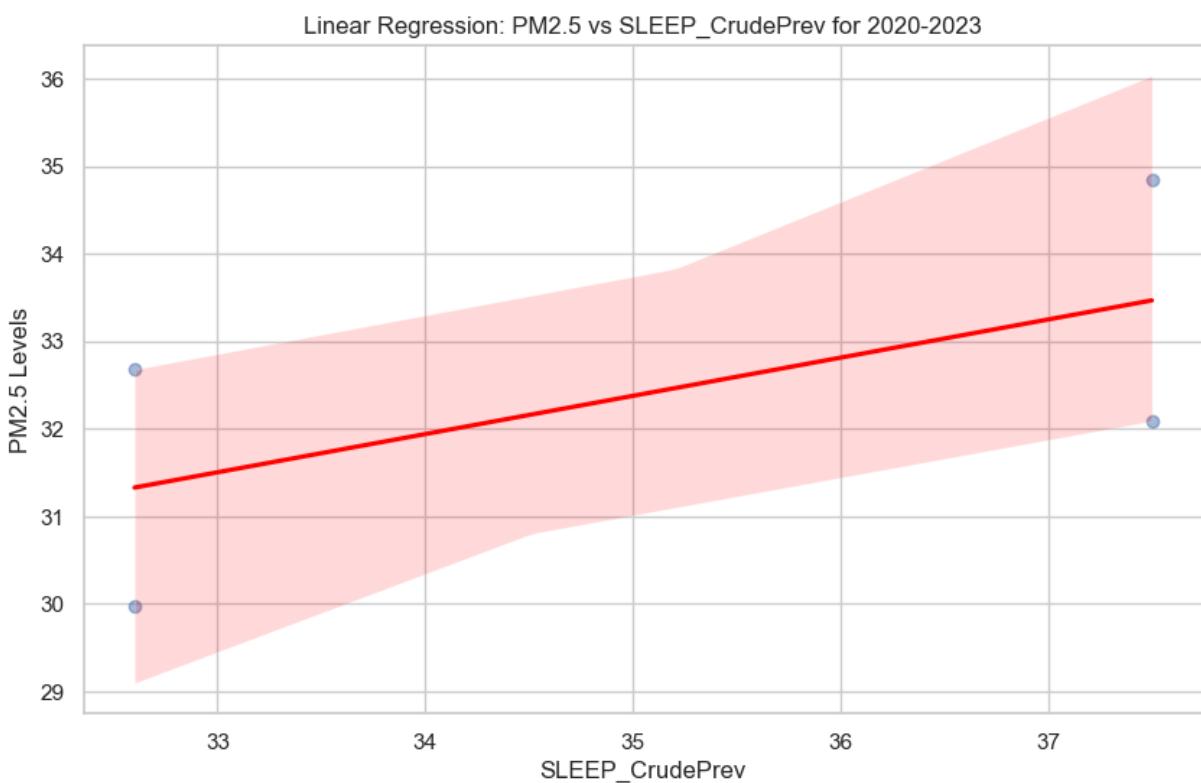
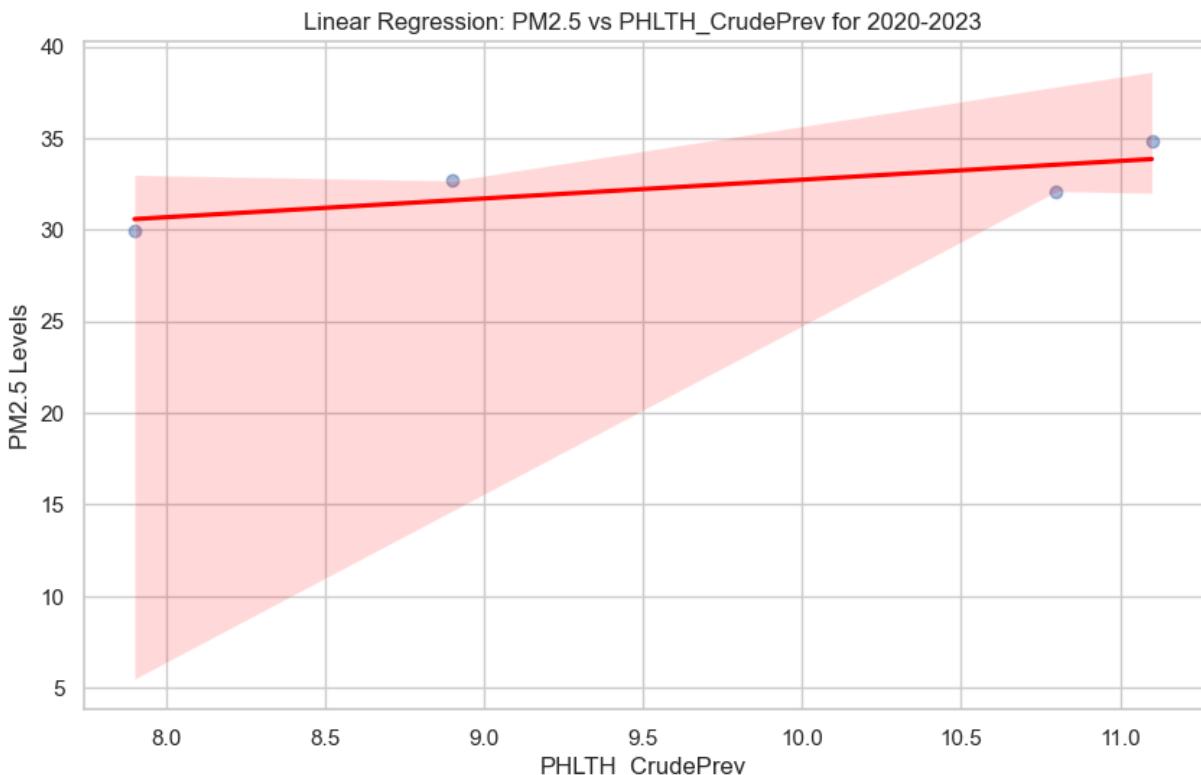


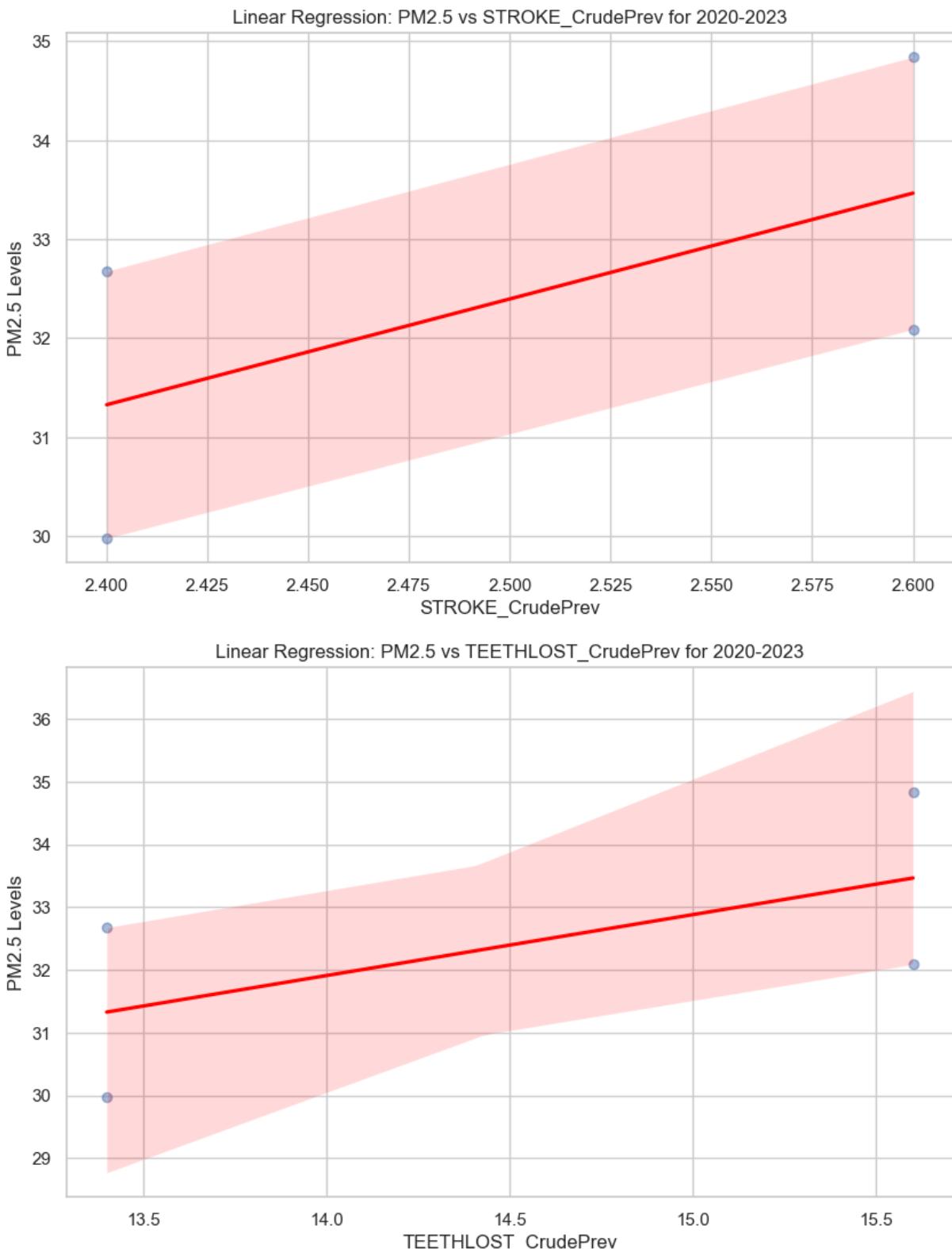
Linear Regression: PM2.5 vs MHLTH_CrudePrev for 2020-2023



Linear Regression: PM2.5 vs OBESITY_CrudePrev for 2020-2023

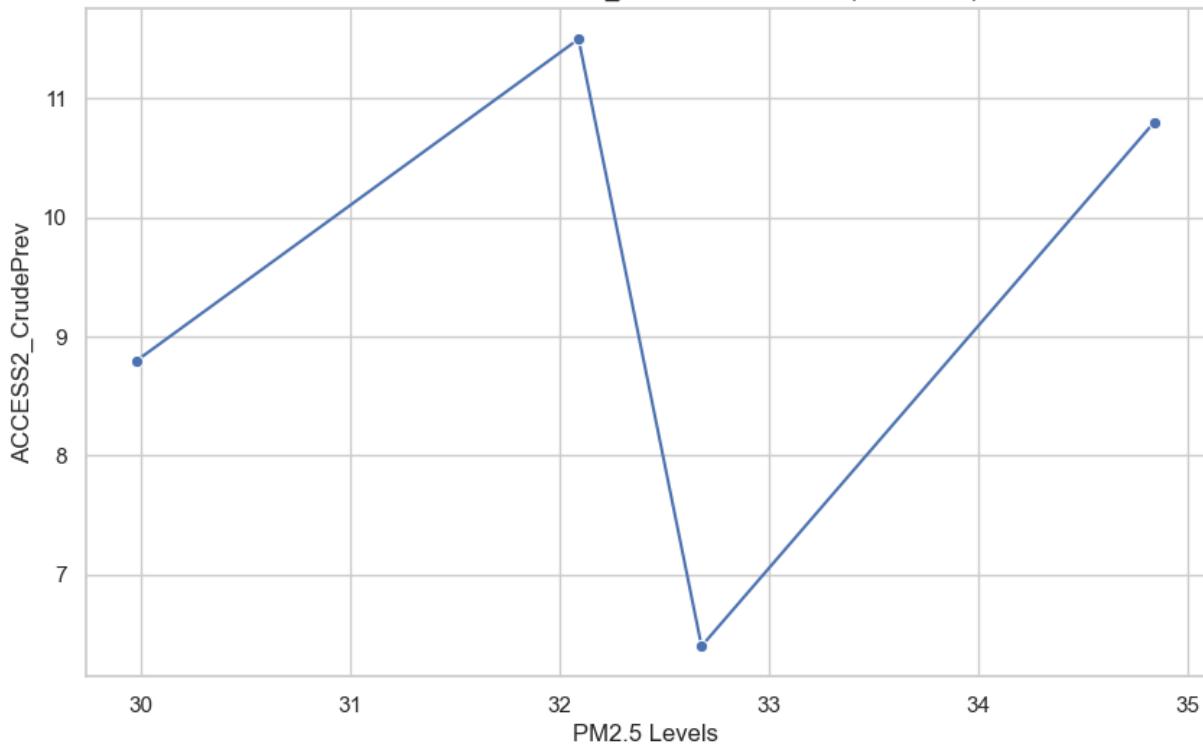




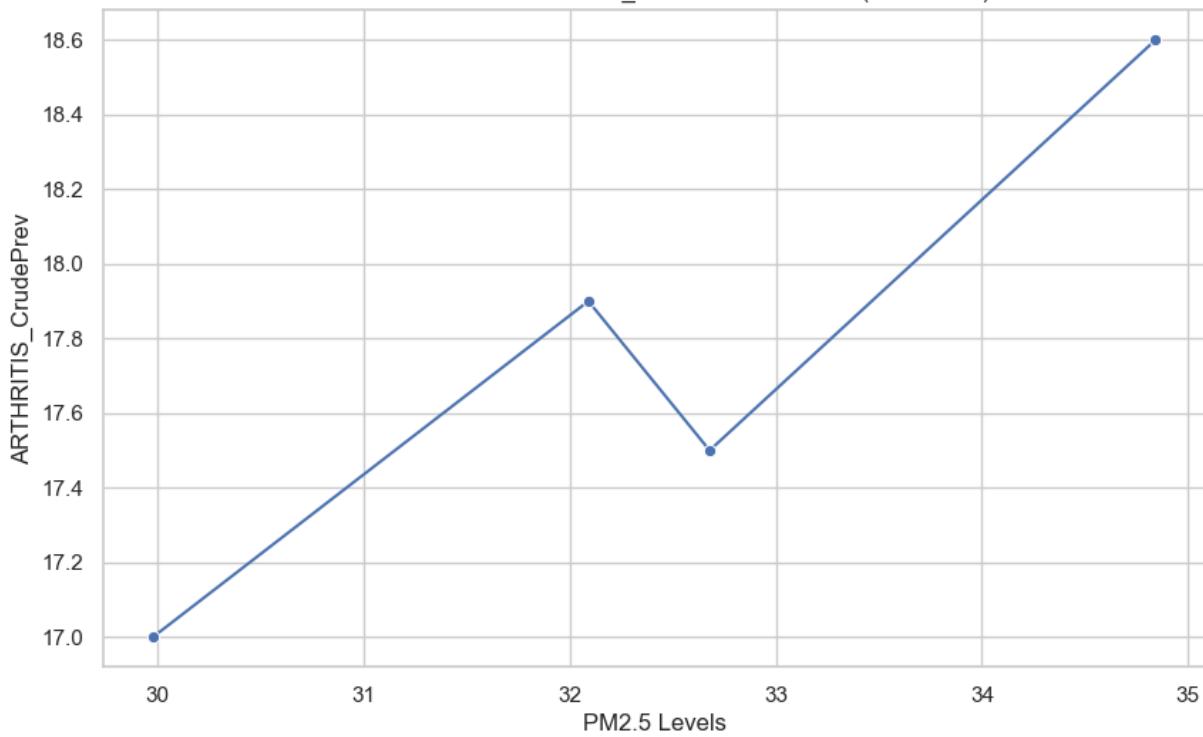


```
In [18]: for column in crude_columns_boston:
    if column != 'pm25':
        plt.figure(figsize=(10, 6))
        sns.lineplot(x='pm25', y=column, data=boston_df, marker='o')
        plt.title(f'PM2.5 Levels vs {column} in Boston (2020-2023)')
        plt.xlabel('PM2.5 Levels')
        plt.ylabel(column)
        plt.show()
```

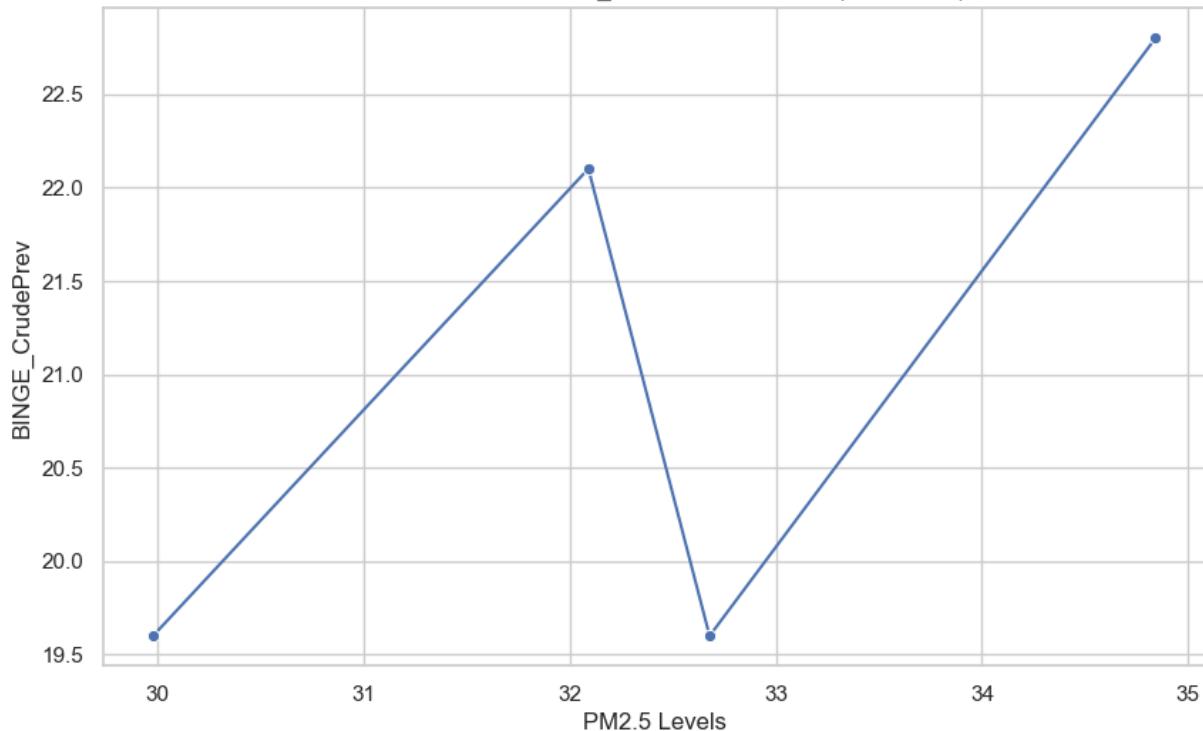
PM2.5 Levels vs ACCESS2_CrudePrev in Boston (2020-2023)



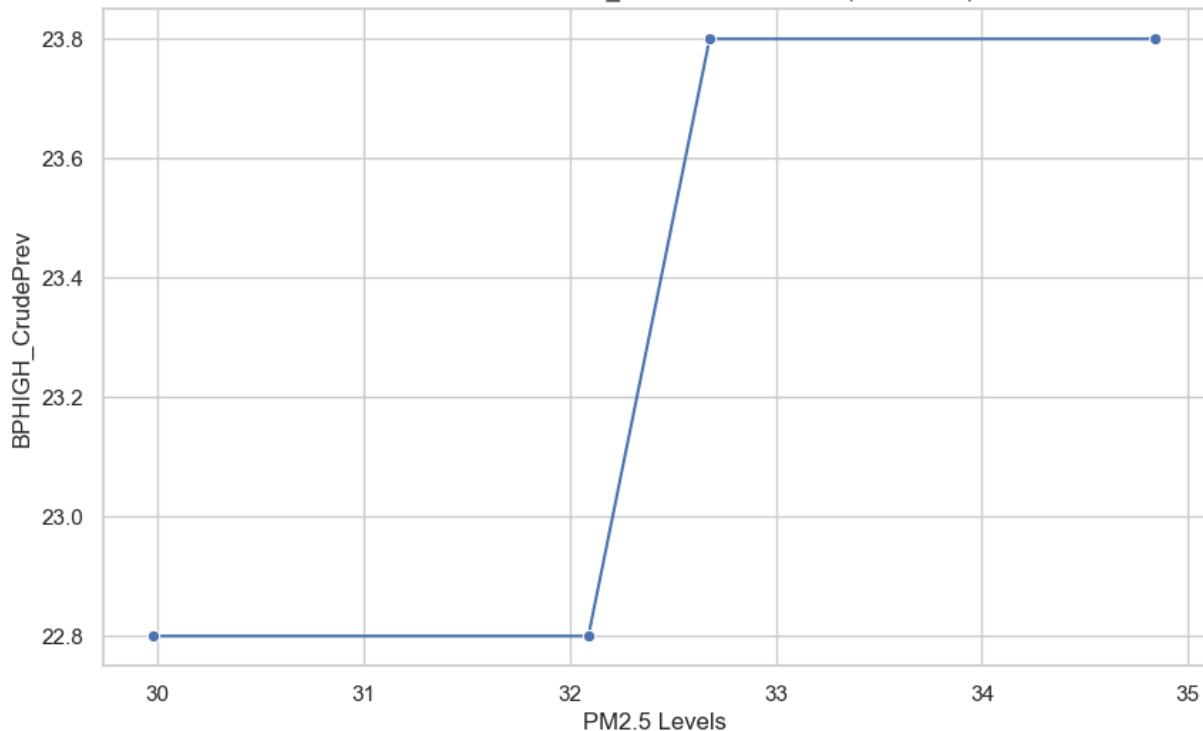
PM2.5 Levels vs ARTHRITIS_CrudePrev in Boston (2020-2023)



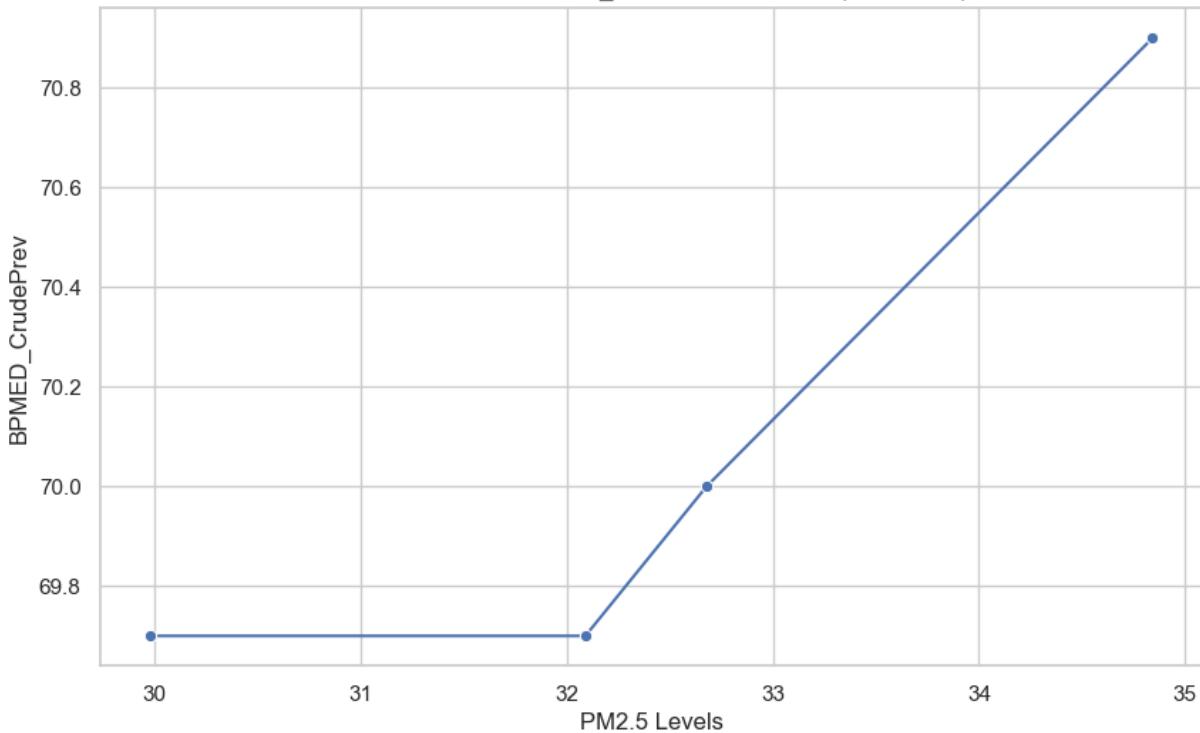
PM2.5 Levels vs BINGE_CrudePrev in Boston (2020-2023)



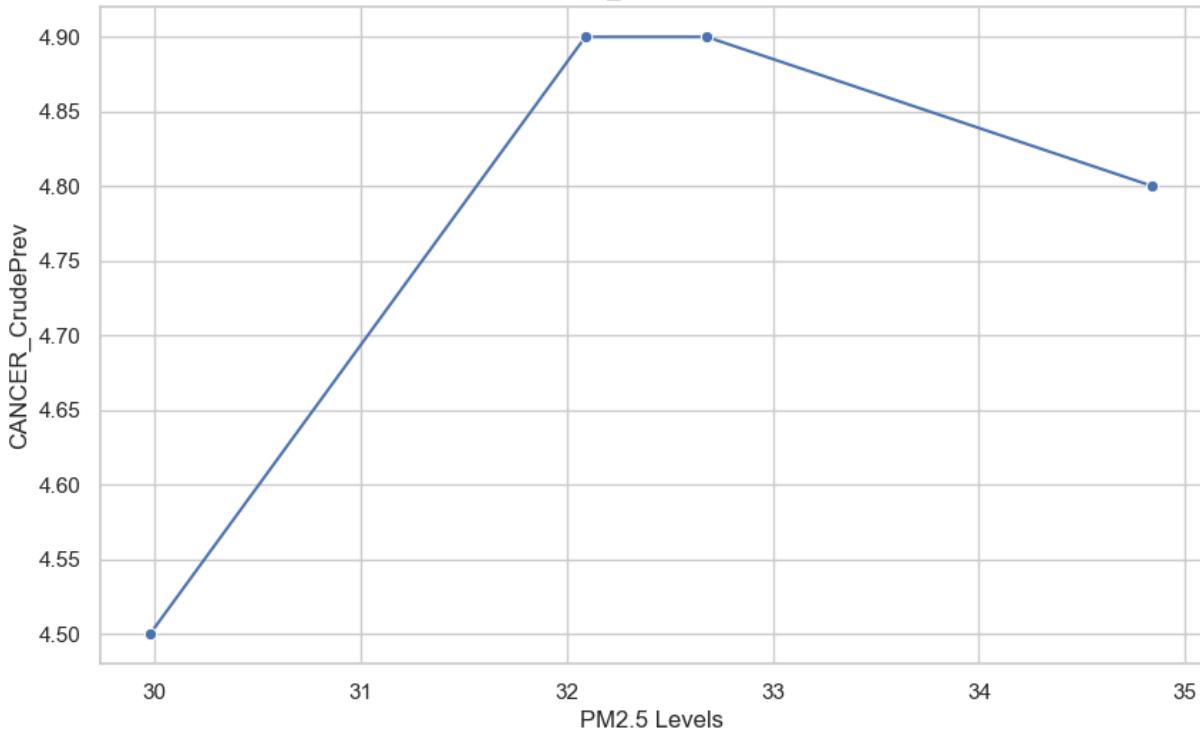
PM2.5 Levels vs BPHIGH_CrudePrev in Boston (2020-2023)



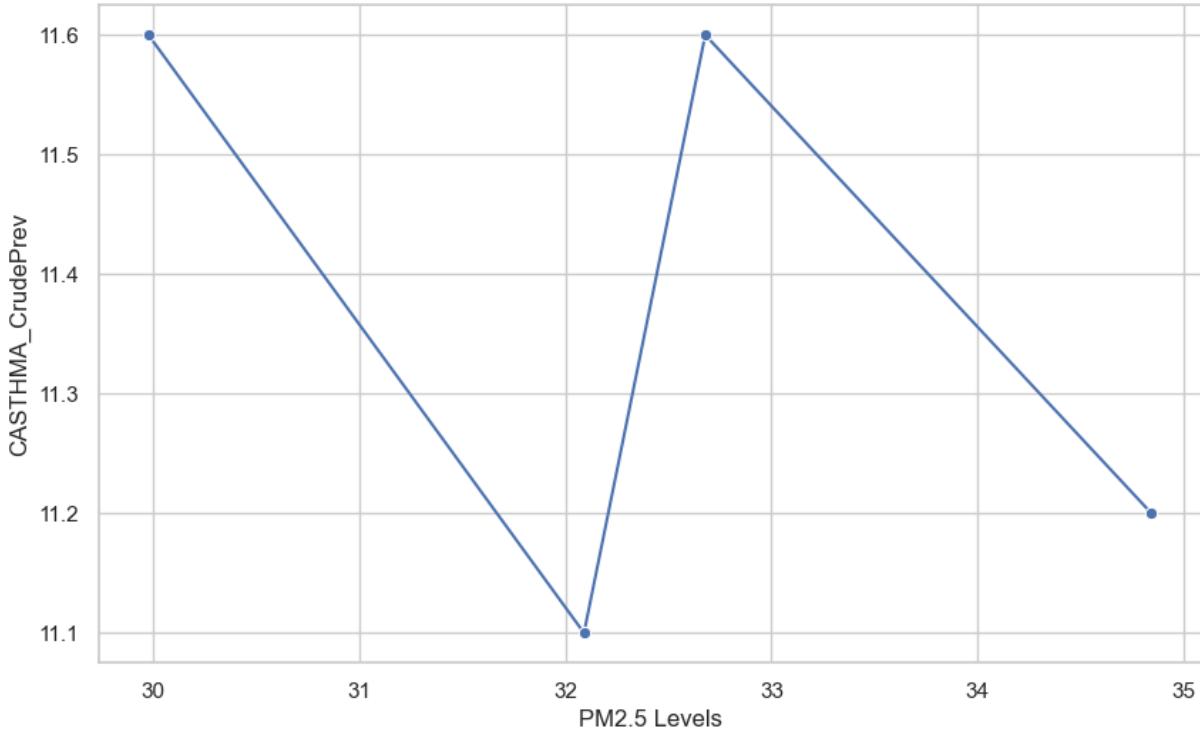
PM2.5 Levels vs BPMED_CrudePrev in Boston (2020-2023)



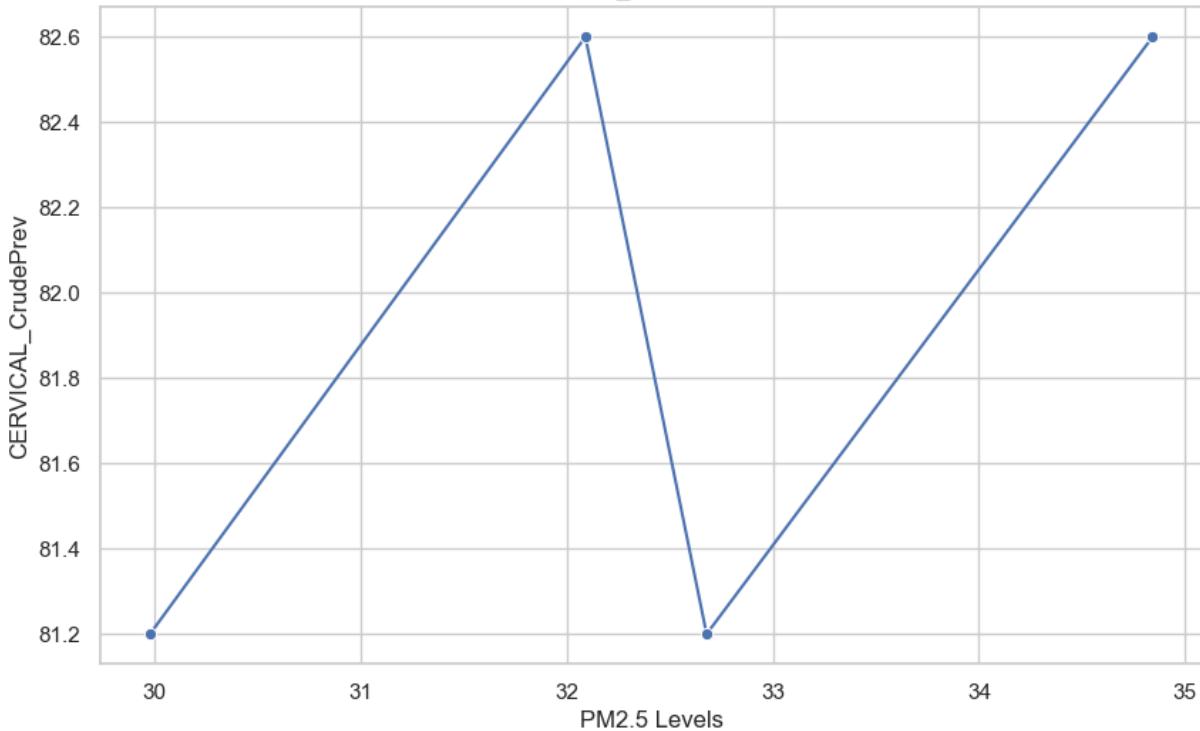
PM2.5 Levels vs CANCER_CrudePrev in Boston (2020-2023)



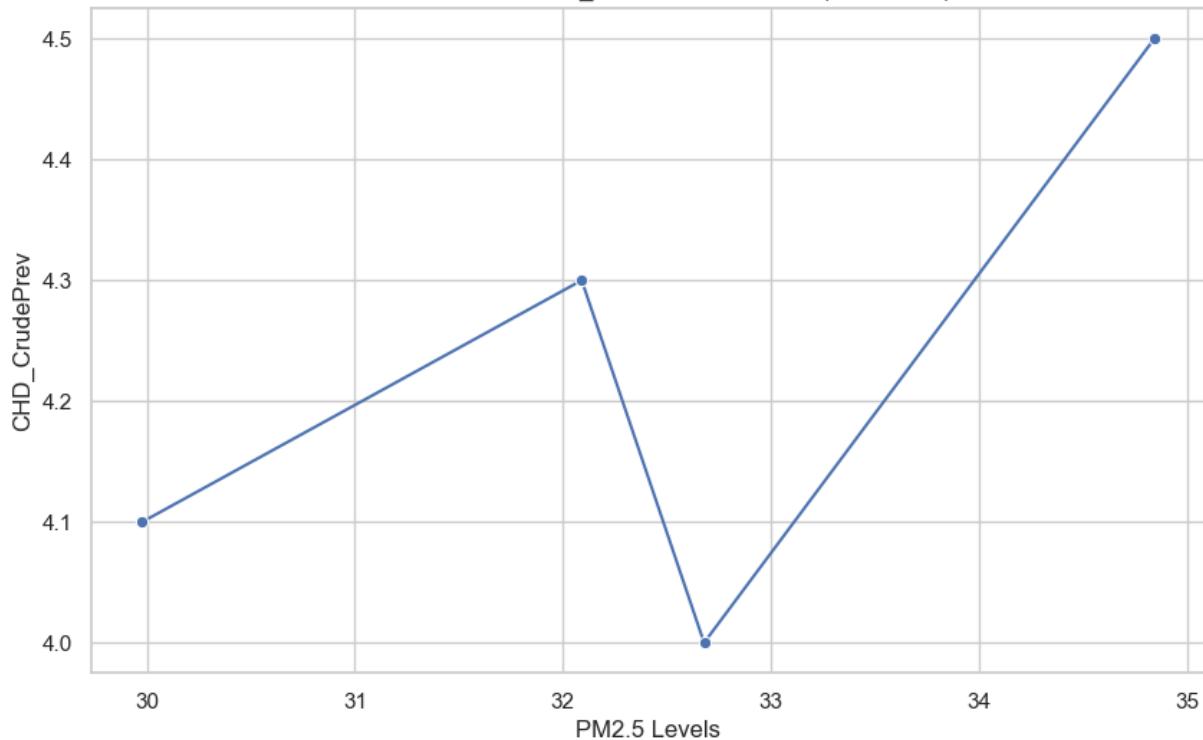
PM2.5 Levels vs CASTHMA_CrudePrev in Boston (2020-2023)



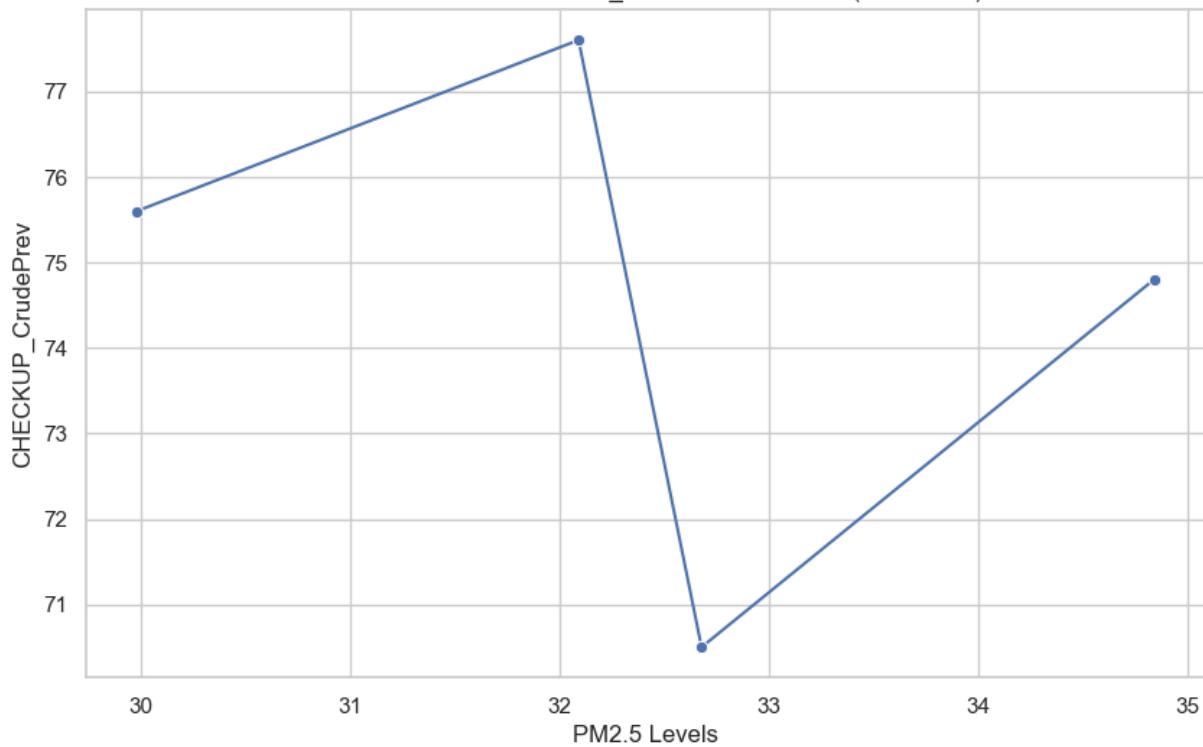
PM2.5 Levels vs CERVICAL_CrudePrev in Boston (2020-2023)

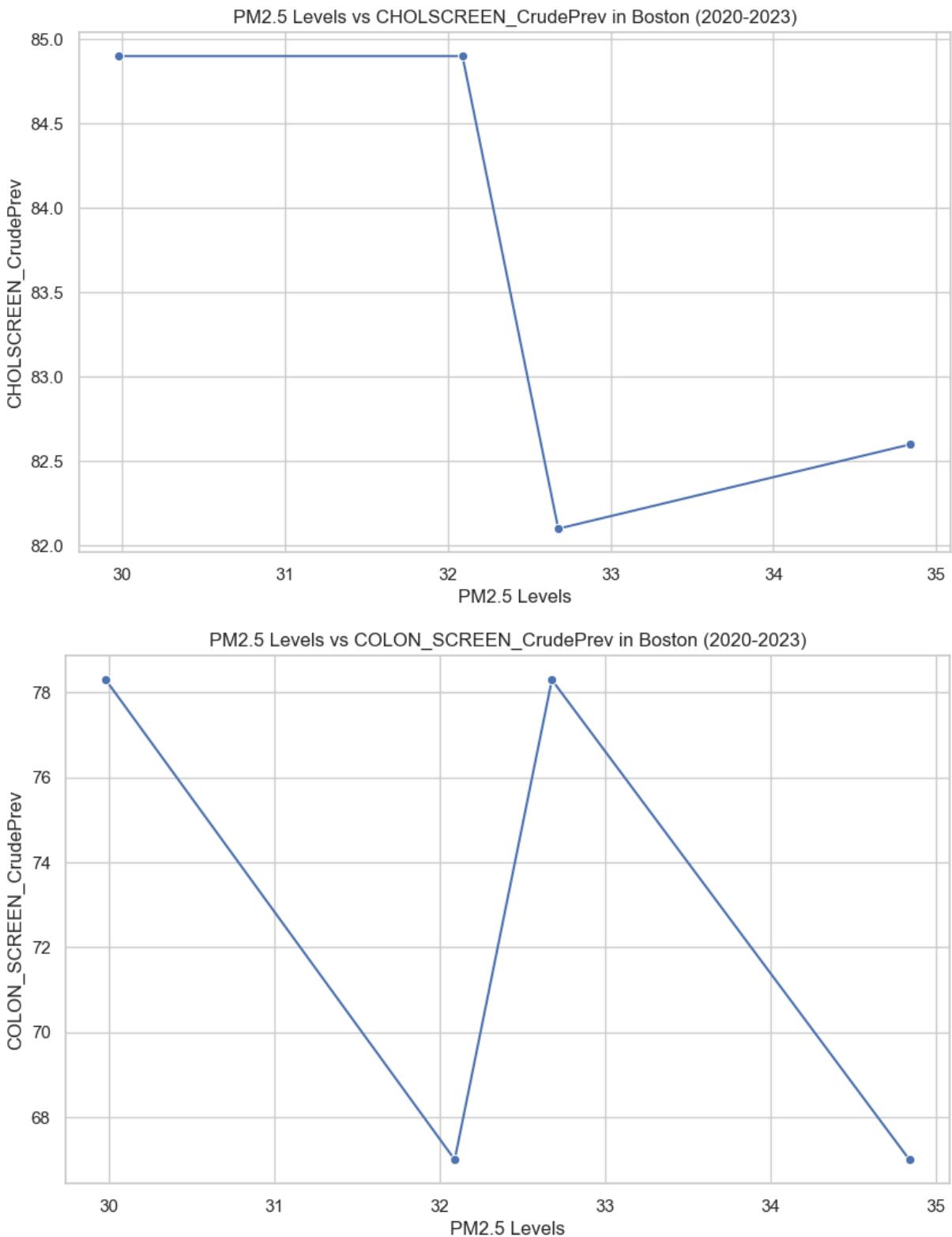


PM2.5 Levels vs CHD_CrudePrev in Boston (2020-2023)

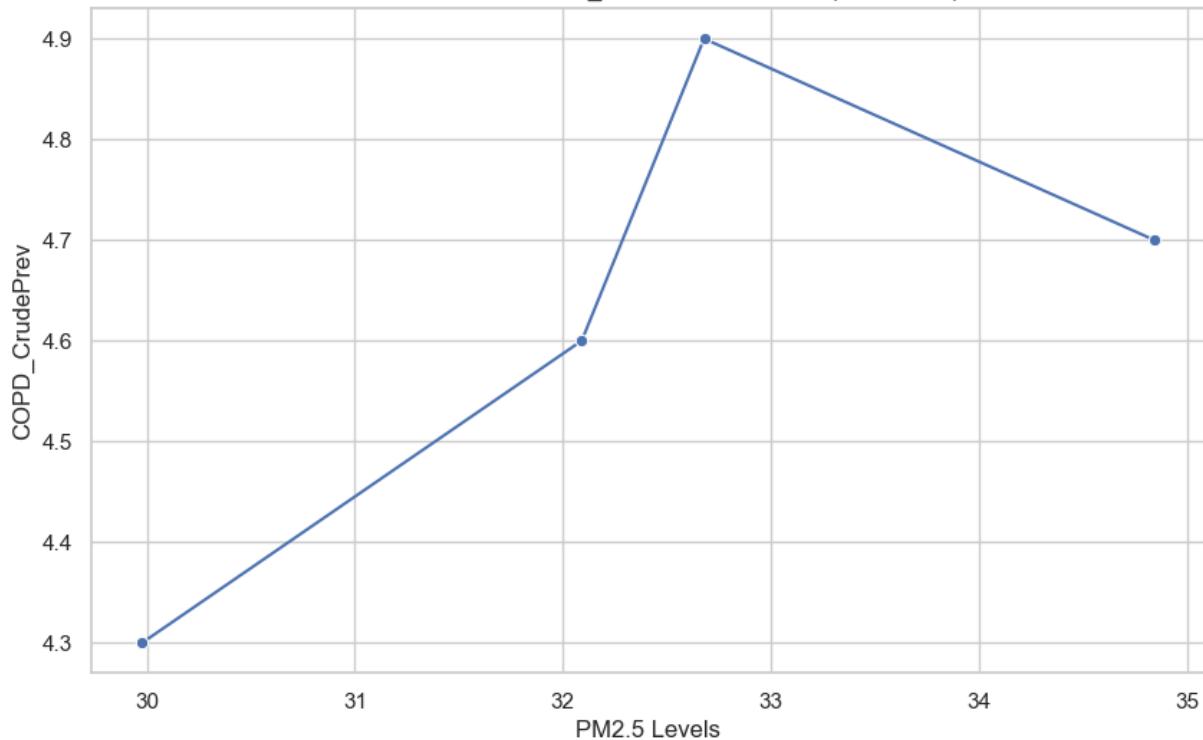


PM2.5 Levels vs CHECKUP_CrudePrev in Boston (2020-2023)

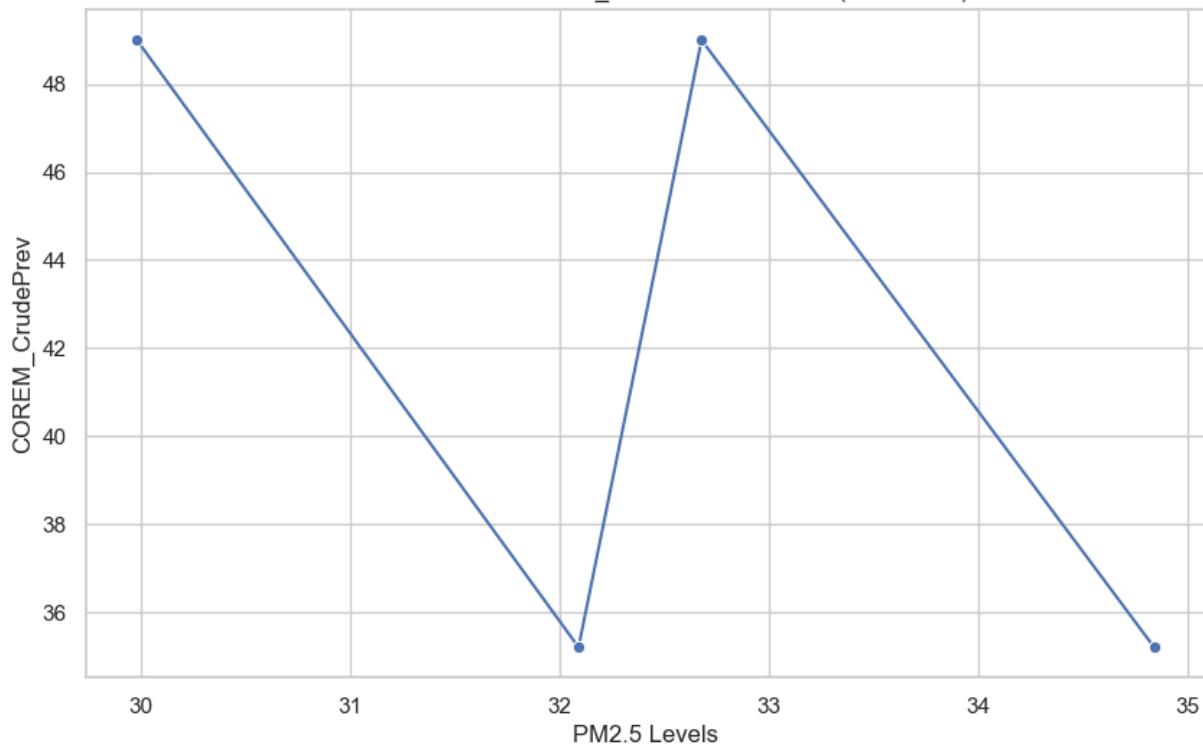




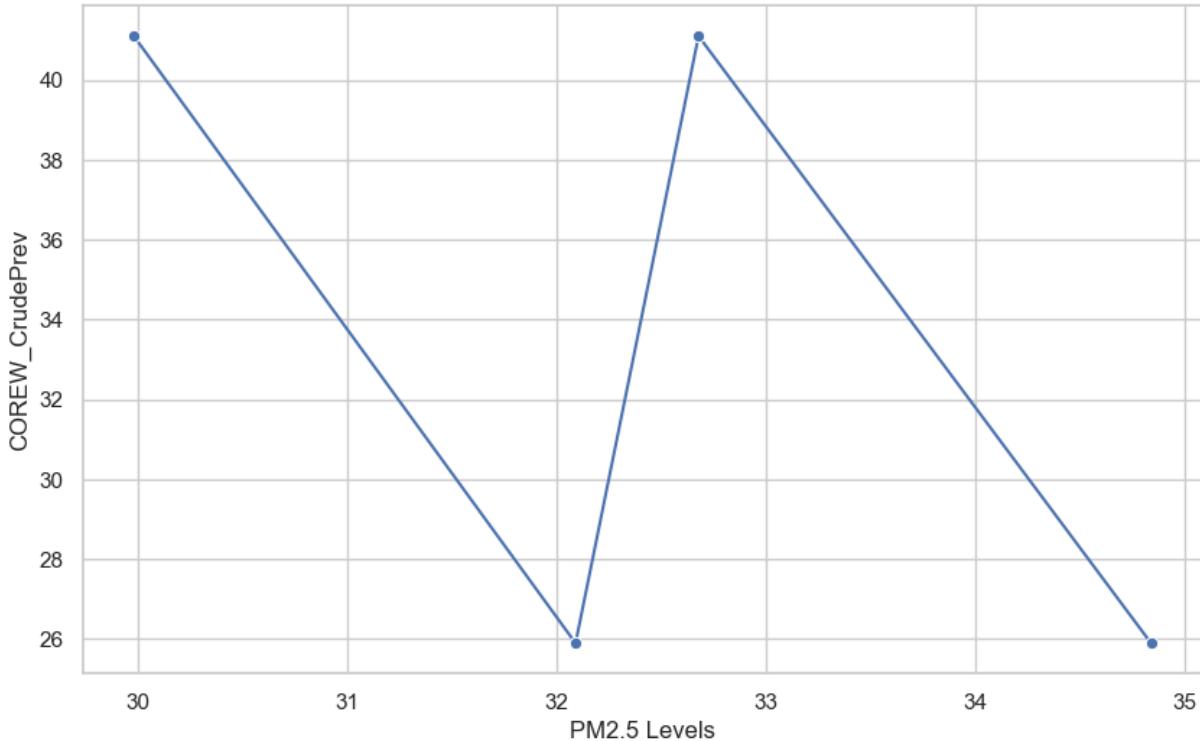
PM2.5 Levels vs COPD_CrudePrev in Boston (2020-2023)



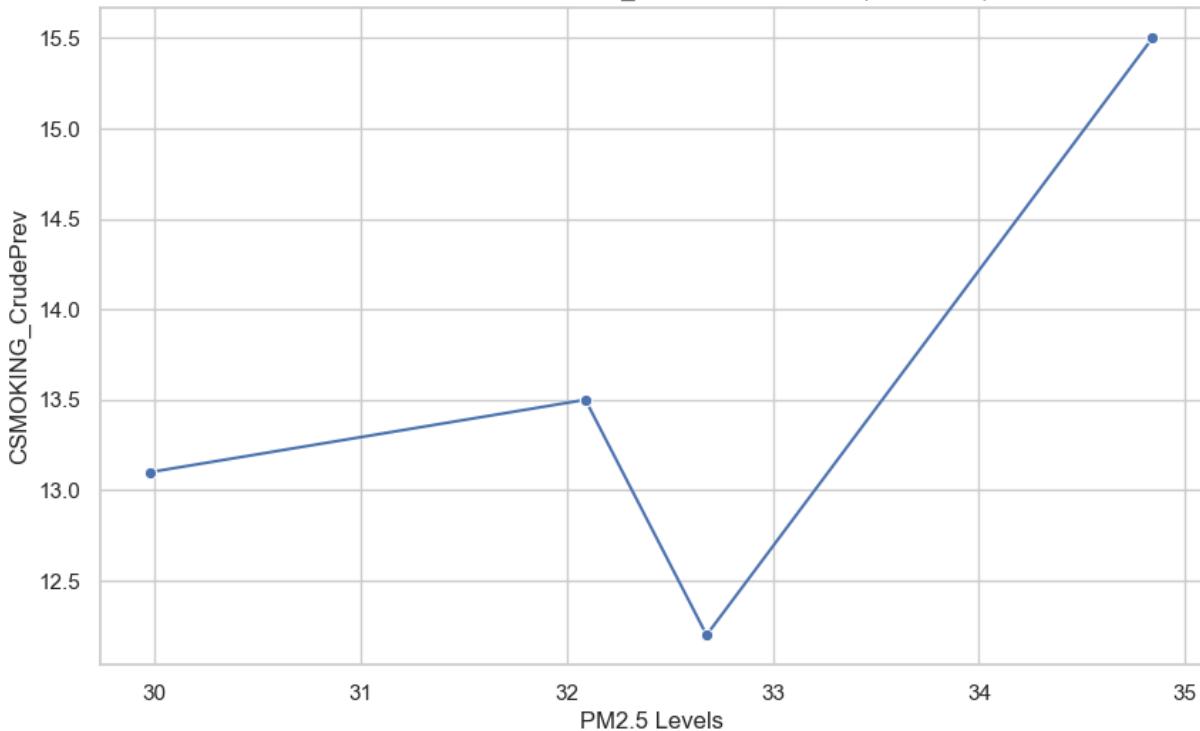
PM2.5 Levels vs COREM_CrudePrev in Boston (2020-2023)



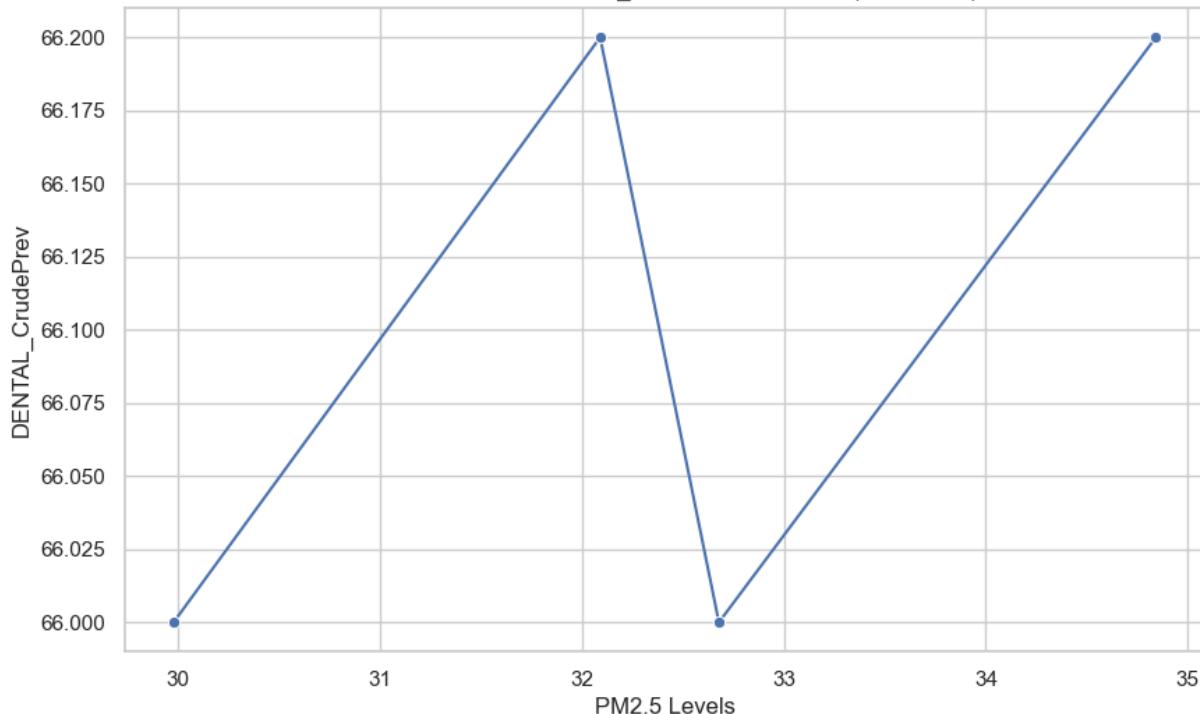
PM2.5 Levels vs COREW_CrudePrev in Boston (2020-2023)



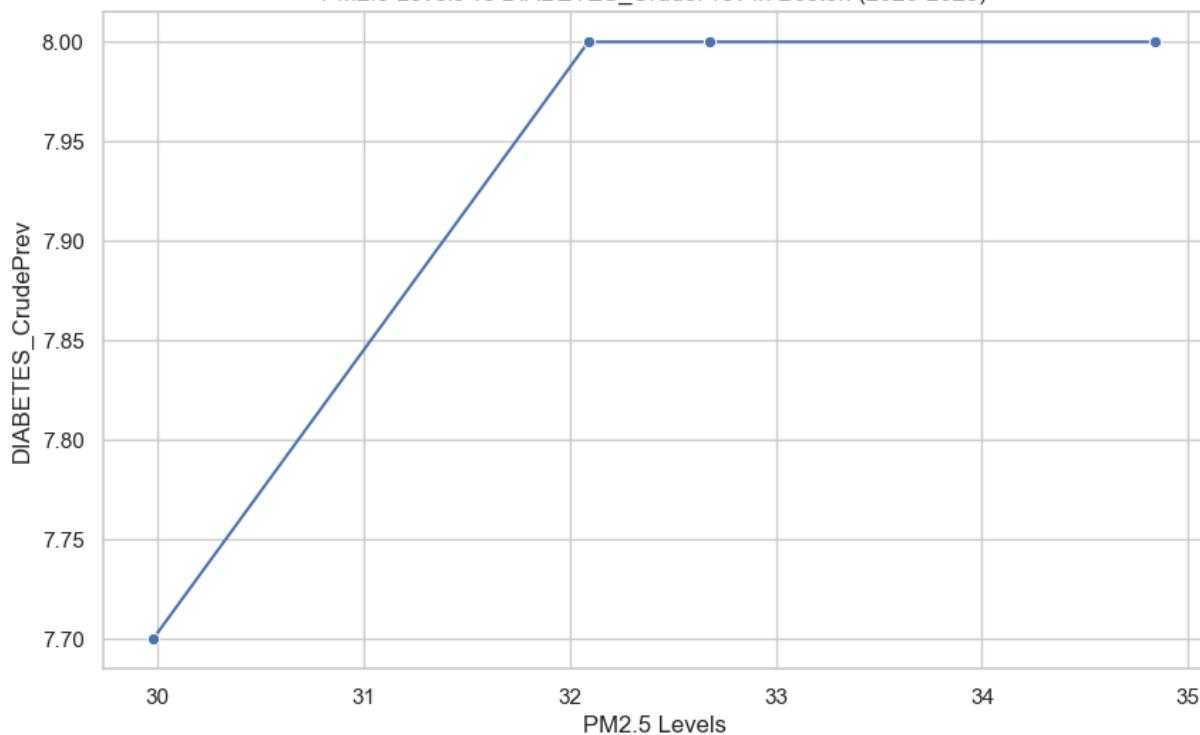
PM2.5 Levels vs CSMOKING_CrudePrev in Boston (2020-2023)



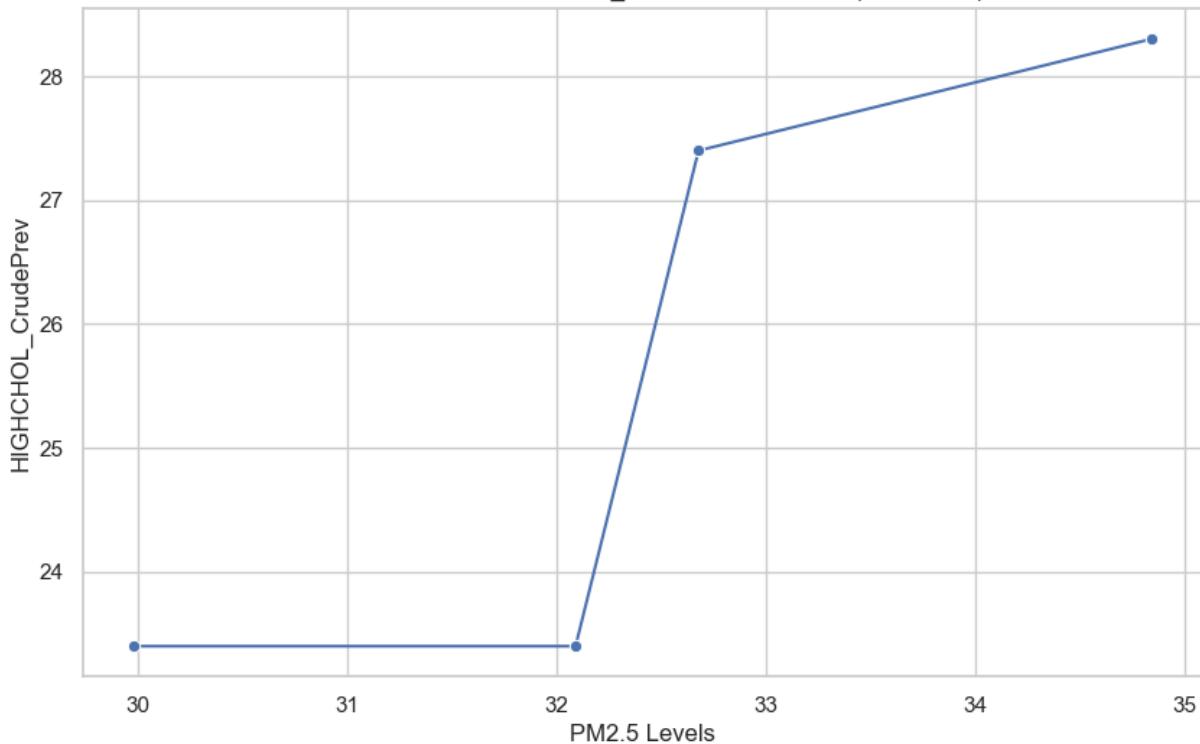
PM2.5 Levels vs DENTAL_CrudePrev in Boston (2020-2023)



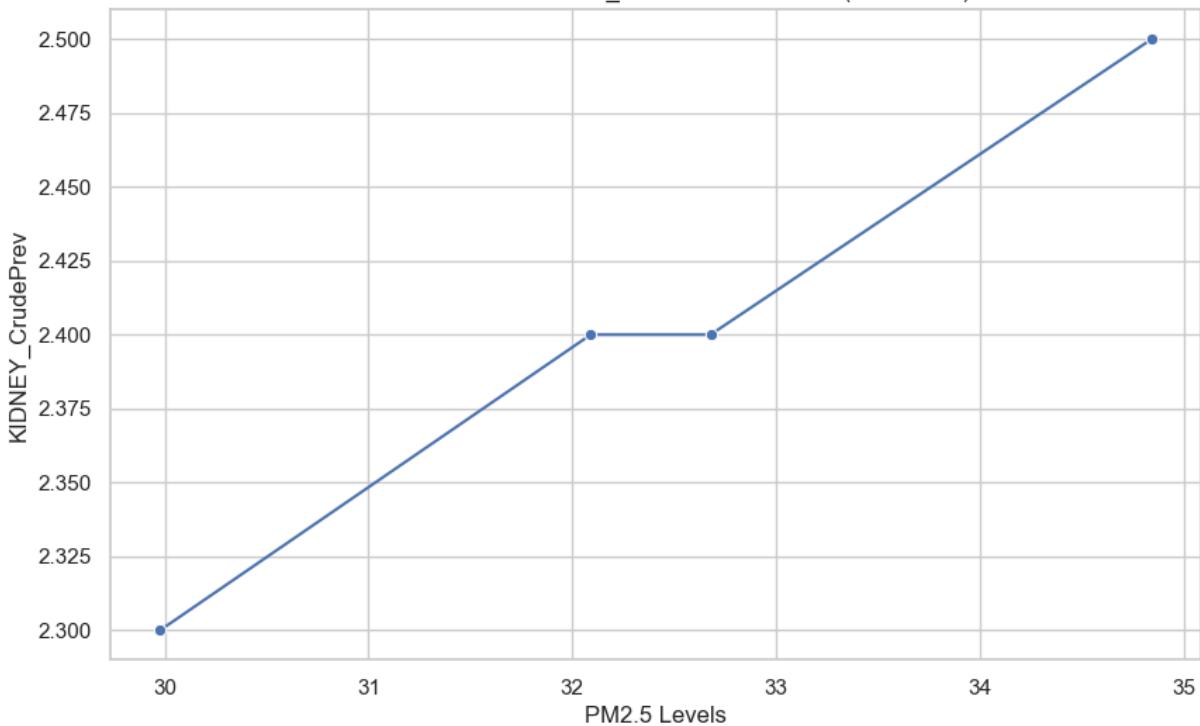
PM2.5 Levels vs DIABETES_CrudePrev in Boston (2020-2023)



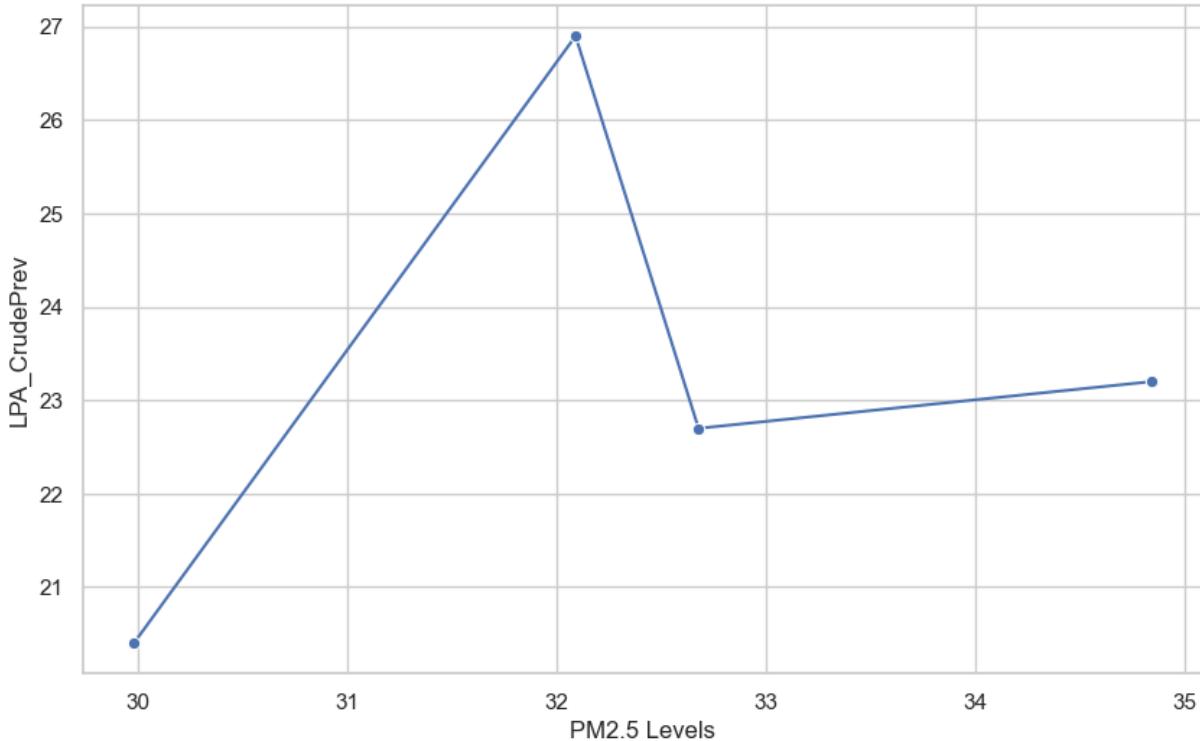
PM2.5 Levels vs HIGHCHOL_CrudePrev in Boston (2020-2023)



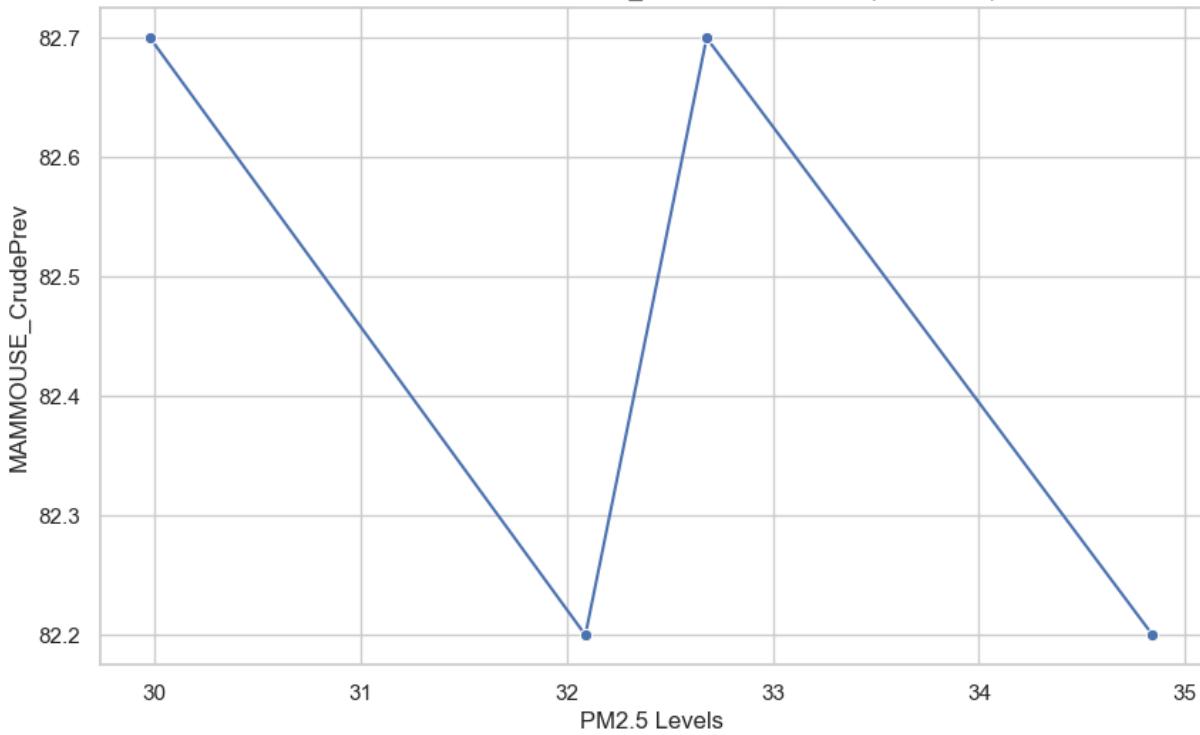
PM2.5 Levels vs KIDNEY_CrudePrev in Boston (2020-2023)



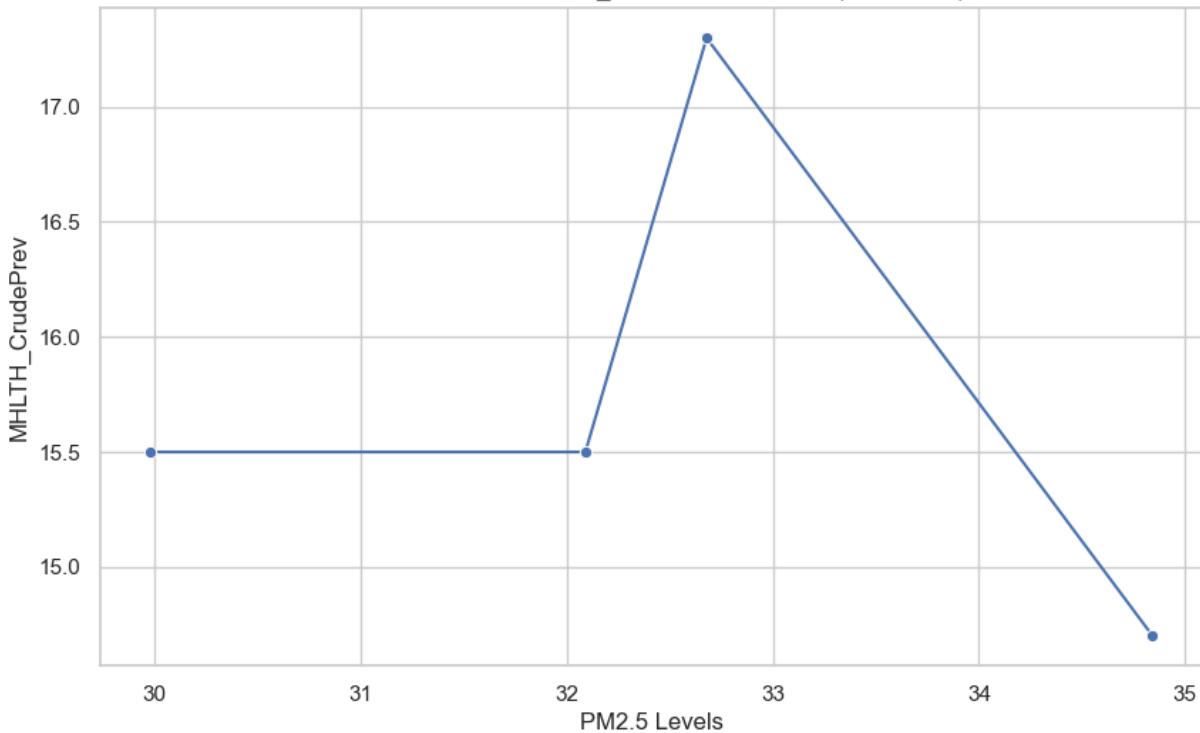
PM2.5 Levels vs LPA_CrudePrev in Boston (2020-2023)



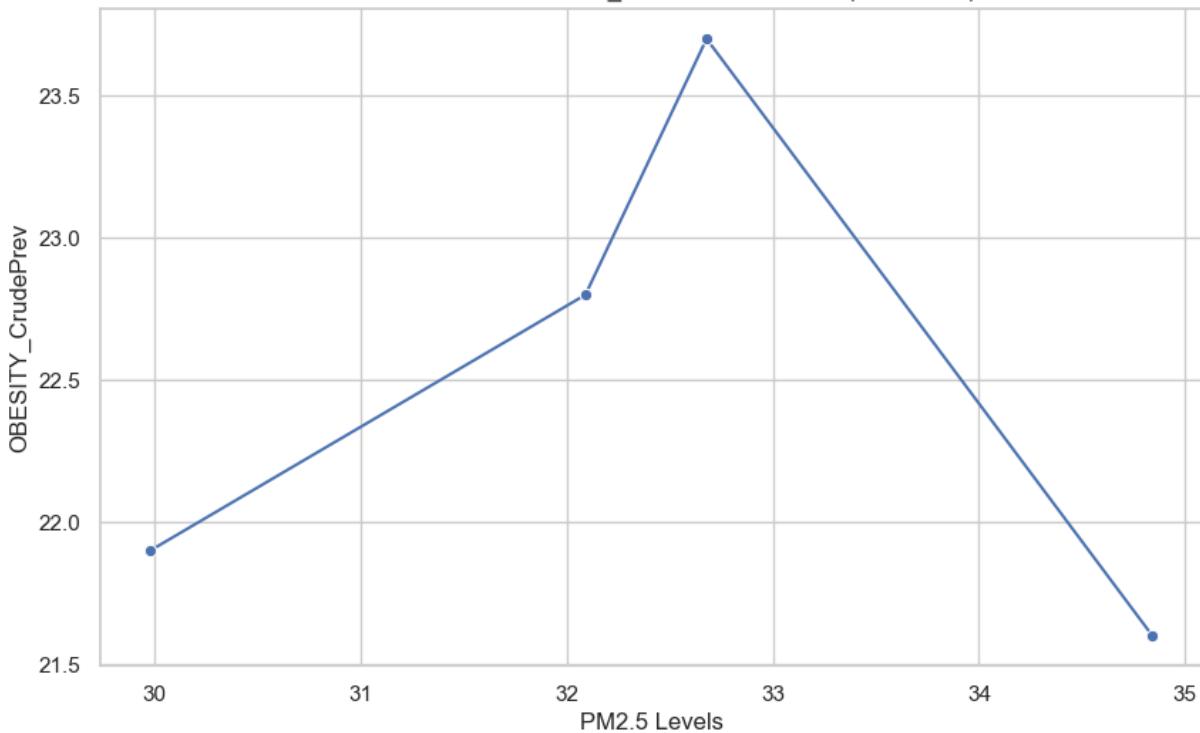
PM2.5 Levels vs MAMMOUSE_CrudePrev in Boston (2020-2023)



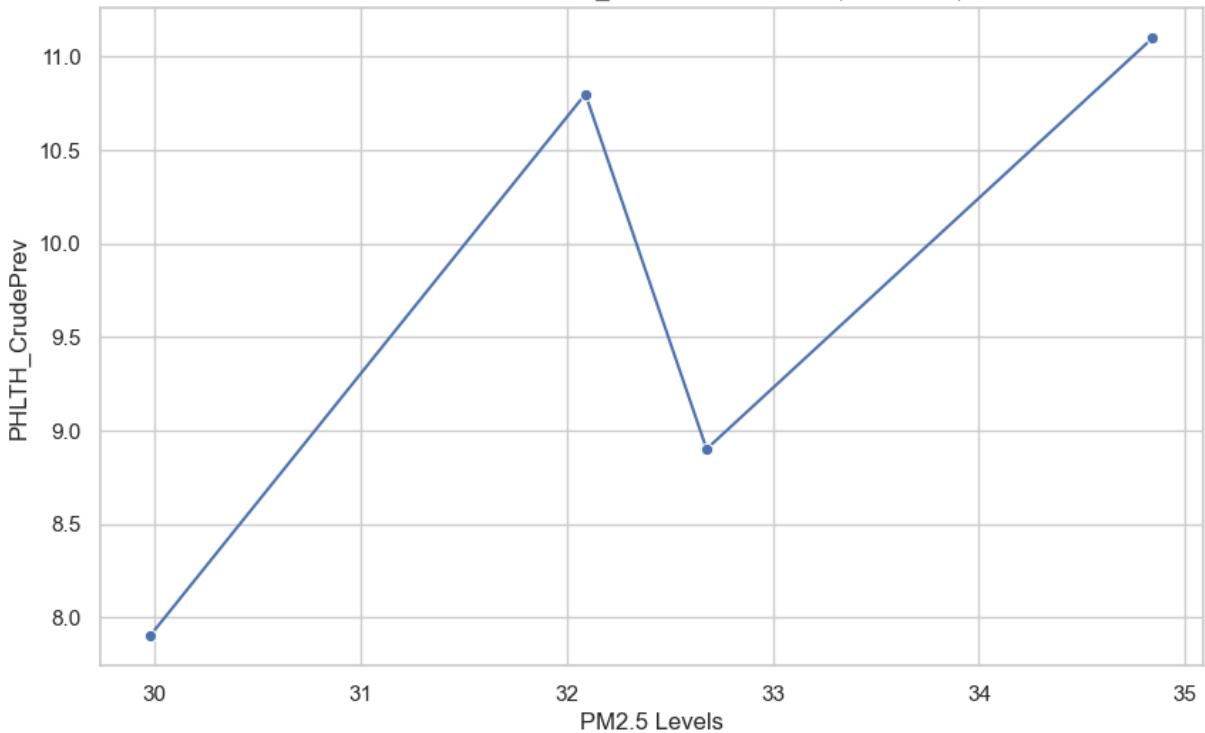
PM2.5 Levels vs MHLTH_CrudePrev in Boston (2020-2023)



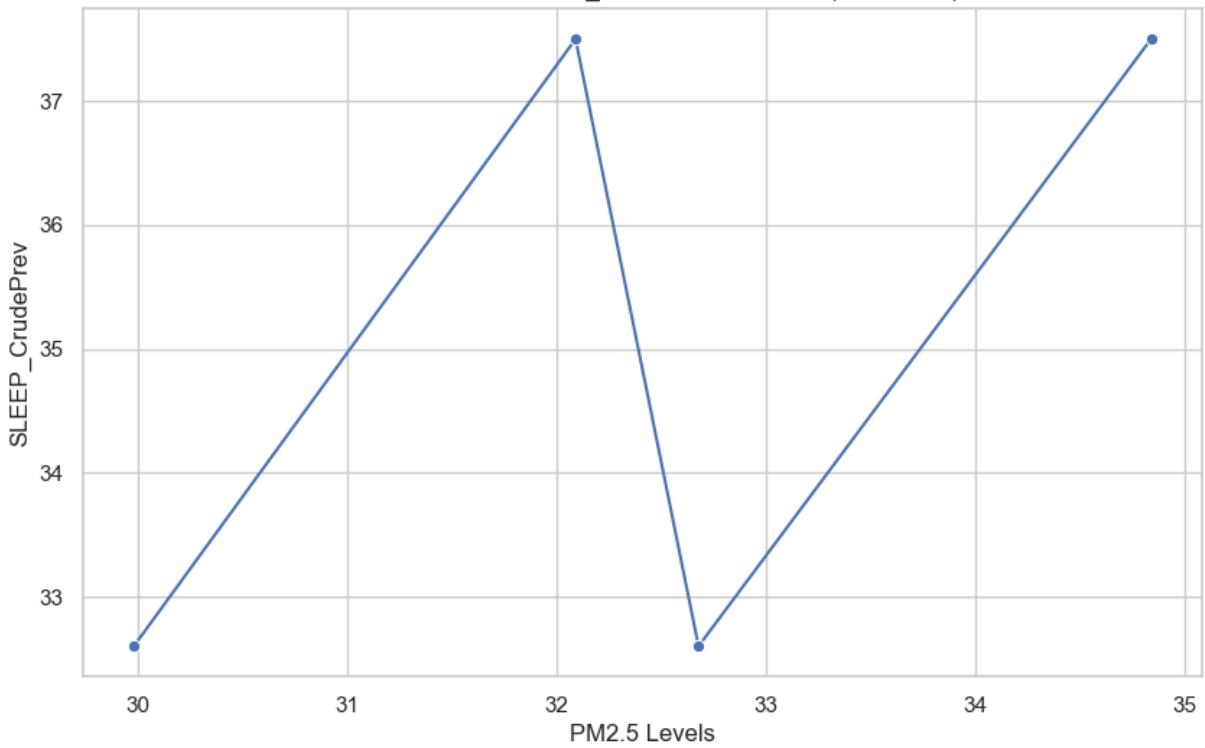
PM2.5 Levels vs OBESITY_CrudePrev in Boston (2020-2023)



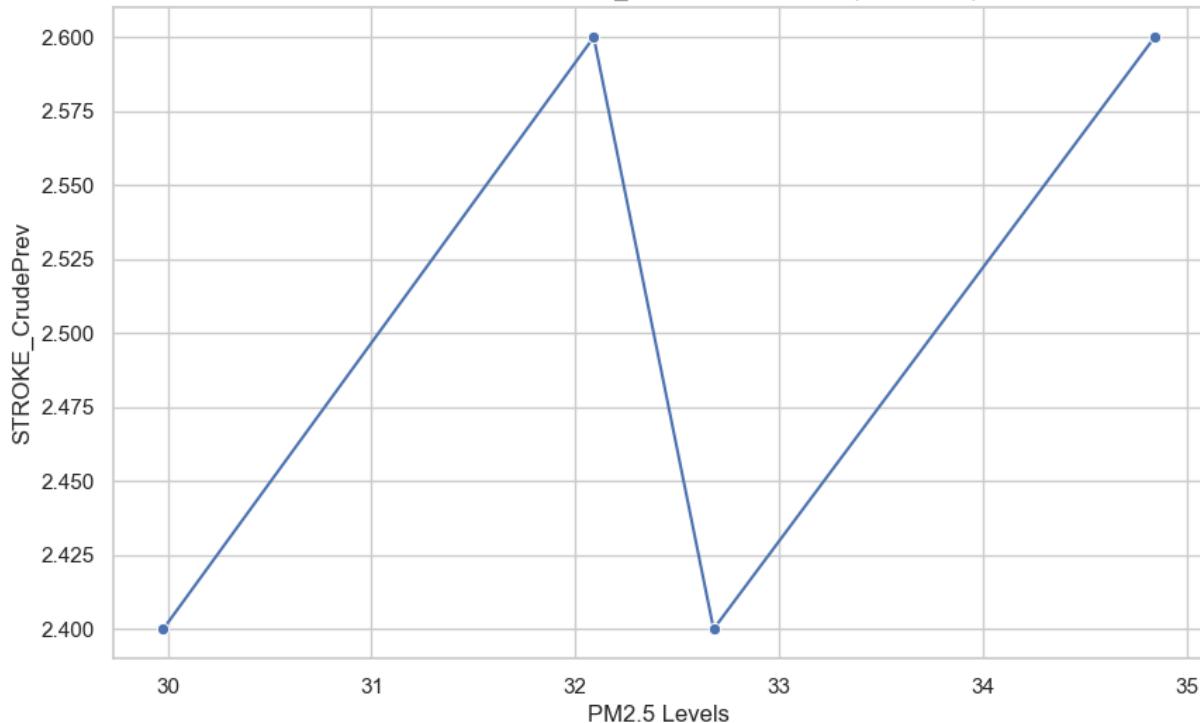
PM2.5 Levels vs PHLTH_CrudePrev in Boston (2020-2023)



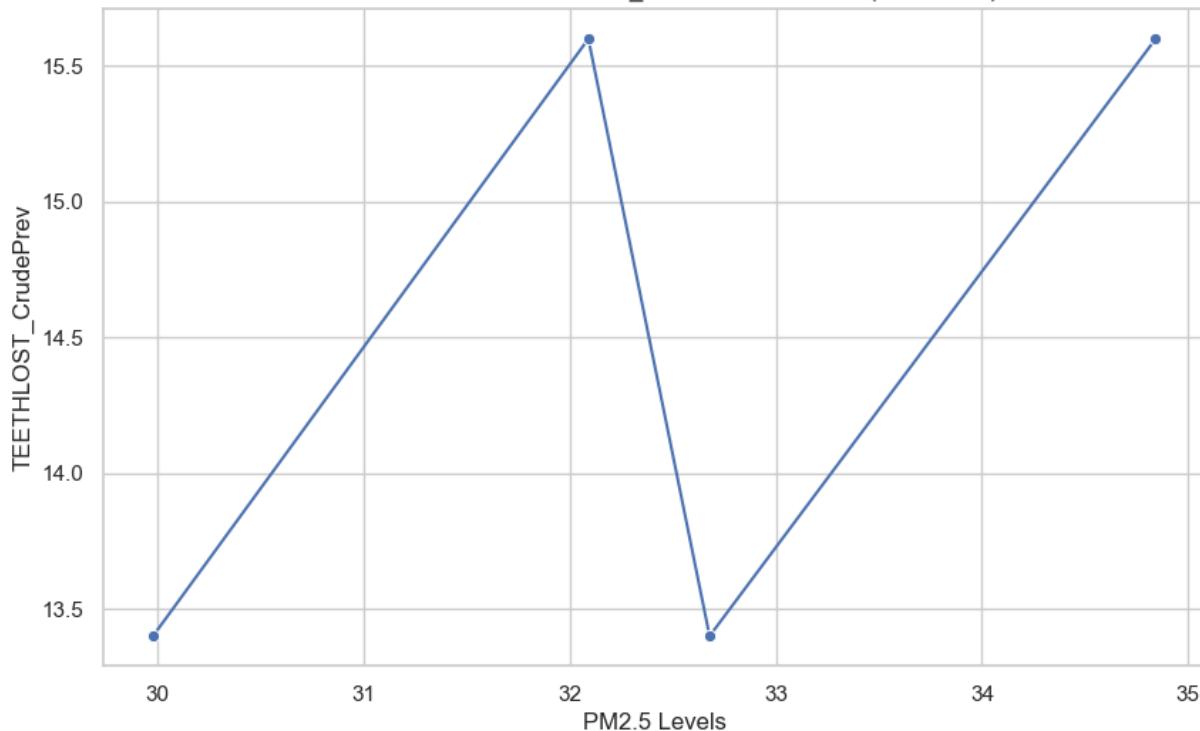
PM2.5 Levels vs SLEEP_CrudePrev in Boston (2020-2023)



PM2.5 Levels vs STROKE_CrudePrev in Boston (2020-2023)



PM2.5 Levels vs TEETHLOST_CrudePrev in Boston (2020-2023)



```
In [19]: non_disease_columns = {'StateAbbr', 'StateDesc', 'PlaceName', 'PlaceFIPS', 'Disease'}
disease_columns = set(boston_df.columns) - non_disease_columns

melted_data = pd.melt(boston_df, id_vars=['Year', 'pm25'], value_vars=disease_columns,
                      var_name='Disease', value_name='CrudePrev')

plt.figure(figsize=(20,10))
sns.lineplot(x='Year', y='CrudePrev', hue='Disease', data=melted_data)
plt.title('Yearly Changes in Crude Rates for Boston (2020-2023)')
```

```

plt.xlabel('Year')
plt.ylabel('Crude Prevalence Rate')
plt.legend(title='Disease', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.grid(True)
plt.show()

```

