

Data Collection and Preprocessing Phase

Date	11 July 2024
Team ID	SWTID1720359900
Project Title	Machine Learning Approach for Predicting The Price Of Natural Gas
Maximum Marks	6 Marks

Data Exploration and Preprocessing

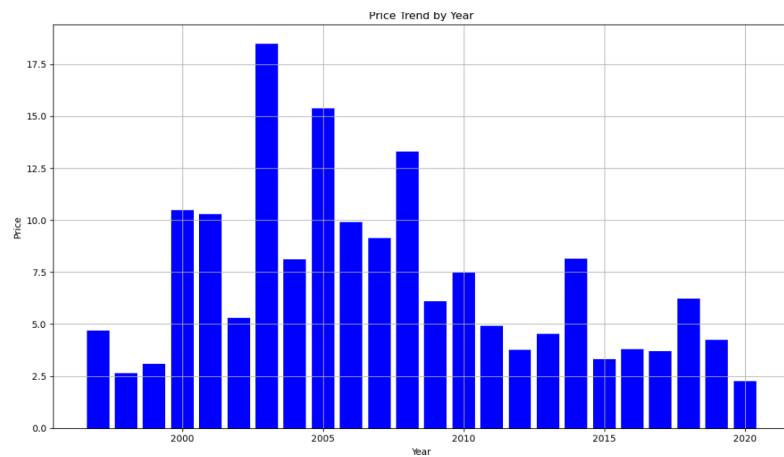
Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

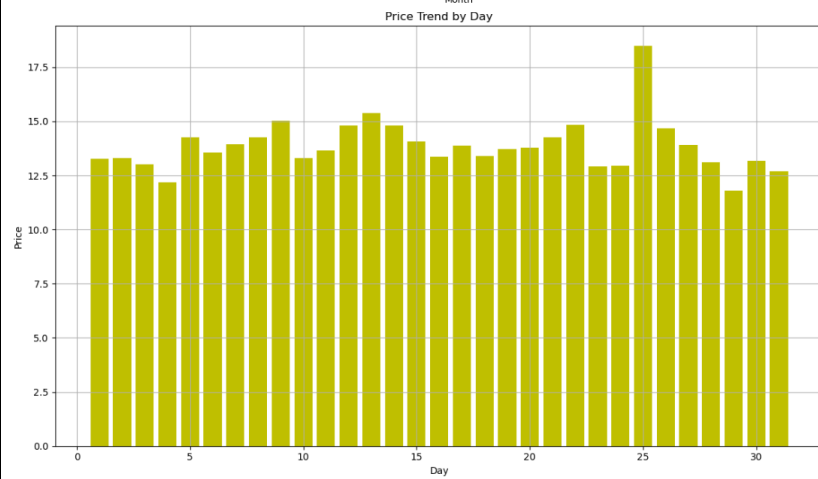
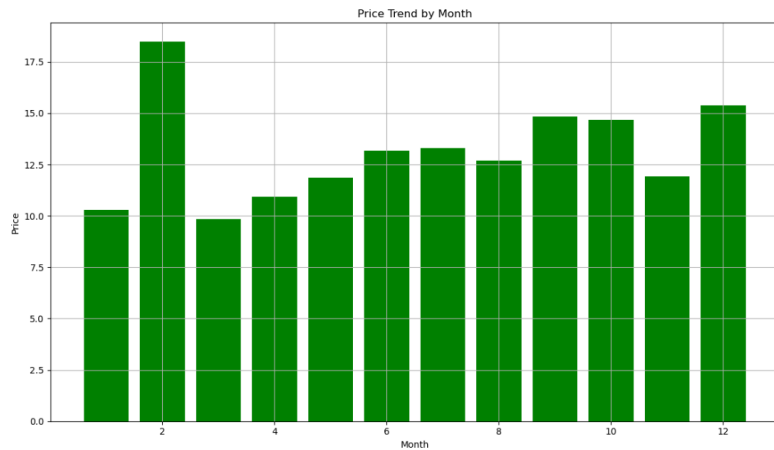
Section	Description																		
Data Overview	<pre>[6]:</pre>																		
	<table><tr><th></th><th>Date</th><th>Price</th></tr><tr><td>0</td><td>1997-01-07</td><td>3.82</td></tr><tr><td>1</td><td>1997-01-08</td><td>3.80</td></tr><tr><td>2</td><td>1997-01-09</td><td>3.61</td></tr><tr><td>3</td><td>1997-01-10</td><td>3.92</td></tr><tr><td>4</td><td>1997-01-13</td><td>4.00</td></tr></table>		Date	Price	0	1997-01-07	3.82	1	1997-01-08	3.80	2	1997-01-09	3.61	3	1997-01-10	3.92	4	1997-01-13	4.00
		Date	Price																
	0	1997-01-07	3.82																
	1	1997-01-08	3.80																
	2	1997-01-09	3.61																
	3	1997-01-10	3.92																
4	1997-01-13	4.00																	
<pre>[5938 rows x 2 columns]></pre>																			

Univariate Analysis

Price	
count	5937.00000
mean	4.18923
std	2.19121
min	1.05000
25%	2.66000
50%	3.54000
75%	5.24000
max	18.48000

Bivariate Analysis





Data Preprocessing Code Screenshots

Loading Data

```
df.head()
```

	Date	Price
0	1997-01-07	3.82
1	1997-01-08	3.80
2	1997-01-09	3.61
3	1997-01-10	3.92
4	1997-01-13	4.00

Handling Missing Data	<pre>df = pd.DataFrame(data) missing_data_rows = df[df.isna().any(axis=1)] print(missing_data_rows)</pre> <table><thead><tr><th></th><th>Date</th><th>Price</th></tr></thead><tbody><tr><td>5284</td><td>2018-01-05</td><td>NaN</td></tr></tbody></table> <pre>df_2 = df.dropna()</pre> <pre>missing_data_rows = df_2[df_2.isna().any(axis=1)] print(missing_data_rows)</pre> <p>Empty DataFrame Columns: [Date, Price] Index: []</p>		Date	Price	5284	2018-01-05	NaN																								
	Date	Price																													
5284	2018-01-05	NaN																													
Data Transformation	<pre>[]: df_2['Date'] = pd.to_datetime(df_2['Date'], infer_datetime_format=True) # Extract year, month, and day from the 'Date' column df_2['year'] = df_2['Date'].dt.year df_2['month'] = df_2['Date'].dt.month df_2['day'] = df_2['Date'].dt.day # Drop the original 'Date' column if not needed df_2.drop(columns=['Date'], inplace=True) # Display the updated DataFrame print(df_2.head())</pre> <table><thead><tr><th></th><th>Price</th><th>year</th><th>month</th><th>day</th></tr></thead><tbody><tr><td>0</td><td>3.82</td><td>1997</td><td>1</td><td>7</td></tr><tr><td>1</td><td>3.80</td><td>1997</td><td>1</td><td>8</td></tr><tr><td>2</td><td>3.61</td><td>1997</td><td>1</td><td>9</td></tr><tr><td>3</td><td>3.92</td><td>1997</td><td>1</td><td>10</td></tr><tr><td>4</td><td>4.00</td><td>1997</td><td>1</td><td>13</td></tr></tbody></table>		Price	year	month	day	0	3.82	1997	1	7	1	3.80	1997	1	8	2	3.61	1997	1	9	3	3.92	1997	1	10	4	4.00	1997	1	13
	Price	year	month	day																											
0	3.82	1997	1	7																											
1	3.80	1997	1	8																											
2	3.61	1997	1	9																											
3	3.92	1997	1	10																											
4	4.00	1997	1	13																											
Feature Engineering	Attached the codes in final submission.																														
Save Processed Data	-																														