

```
1

1 from google.colab import drive
2 drive.mount('/content/drive')

Mounted at /content/drive

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np

1 data = pd.read_csv('/content/drive/MyDrive/AirQualityUCI.csv')

1 data.head()
```

	Date	Time	CO(GT)	PT08.S1(CO)	NMHC(GT)	C6H6(GT)	PT08.S2(NMHC)	NOx(GT)
0	3/10/2004	18:00:00	2.6	1360	150	11.9	1046	
1	3/10/2004	19:00:00	2.0	1292	112	9.4	955	
2	3/10/2004	20:00:00	2.2	1402	88	9.0	939	
3	3/10/2004	21:00:00	2.2	1376	80	9.2	948	
4	3/10/2004	22:00:00	1.6	1272	51	6.5	836	

```
1 data.columns=['DATE','TIME','CO_GT','PT08_S1_CO','NMHC_GT','C6H6_GT','PT08_S2_NMHC','NOX_GT','PT08_S3_NOX','NO2_GT','PT08_S4_NO2','PT08_S5_O3','T','RH','AH']

1 data.shape

(9357, 15)

1 data.head()
```

	DATE	TIME	CO_GT	PT08_S1_CO	NMHC_GT	C6H6_GT	PT08_S2_NMHC	NOX_GT	P
0	3/10/2004	18:00:00	2.6	1360	150	11.9	1046	166	
1	3/10/2004	19:00:00	2.0	1292	112	9.4	955	103	
2	3/10/2004	20:00:00	2.2	1402	88	9.0	939	131	
3	3/10/2004	21:00:00	2.2	1376	80	9.2	948	172	
4	3/10/2004	22:00:00	1.6	1272	51	6.5	836	131	

```
1 data.replace(to_replace= -200, value= np.NaN, inplace= True)

1 data.isna().sum()

DATE                0
TIME                0
CO_GT              1683
PT08_S1_CO         366
NMHC_GT           8443
C6H6_GT           366
PT08_S2_NMHC       366
NOX_GT            1639
PT08_S3_NOX        366
NO2_GT            1642
PT08_S4_NO2        366
PT08_S5_O3         366
T                  366
RH                 366
AH                 366
dtype: int64

1 data.drop('NMHC_GT',axis=1,inplace=True)

1 data.head()
```

	DATE	TIME	CO_GT	PT08_S1_CO	C6H6_GT	PT08_S2_NMHC	NOX_GT	PT08_S3_NO
0	3/10/2004	18:00:00	2.6	1360.0	11.9	1046.0	166.0	1056.
1	3/10/2004	19:00:00	2.0	1292.0	9.4	955.0	103.0	1174.
2	3/10/2004	20:00:00	2.2	1402.0	9.0	939.0	131.0	1140.

1 data.isna().sum()

DATE	0
TIME	0
CO_GT	1683
PT08_S1_CO	366
C6H6_GT	366
PT08_S2_NMHC	366
NOX_GT	1639
PT08_S3_NOX	366
NO2_GT	1642
PT08_S4_NO2	366
PT08_S5_O3	366
T	366
RH	366
AH	366
dtype:	int64

1 data.dtypes

DATE	object
TIME	object
CO_GT	float64
PT08_S1_CO	float64
C6H6_GT	float64
PT08_S2_NMHC	float64
NOX_GT	float64
PT08_S3_NOX	float64
NO2_GT	float64
PT08_S4_NO2	float64
PT08_S5_O3	float64
T	float64
RH	float64
AH	float64
dtype:	object

1 data.describe()

	CO_GT	PT08_S1_CO	C6H6_GT	PT08_S2_NMHC	NOX_GT	PT08_S3_NO
count	7674.000000	8991.000000	8991.000000	8991.000000	7718.000000	8991.000000
mean	2.152750	1099.833166	10.083105	939.153376	246.896735	835.493609
std	1.453252	217.080037	7.449820	266.831429	212.979168	256.817320
min	0.100000	647.000000	0.100000	383.000000	2.000000	322.000000
25%	1.100000	937.000000	4.400000	734.500000	98.000000	658.000000
50%	1.800000	1063.000000	8.200000	909.000000	180.000000	806.000000
75%	2.900000	1231.000000	14.000000	1116.000000	326.000000	969.500000
max	11.900000	2040.000000	63.700000	2214.000000	1479.000000	2683.000000

1 data.dropna().shape

(6941, 14)

1 data.dropna(inplace=True)

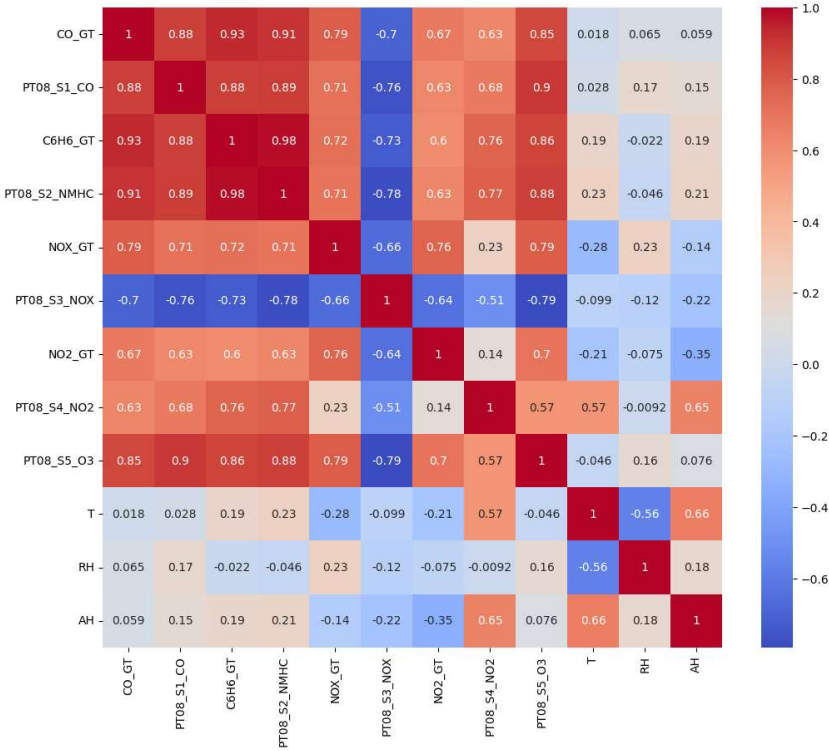
1 data.isnull().sum()

DATE	0
TIME	0
CO_GT	0
PT08_S1_CO	0
C6H6_GT	0
PT08_S2_NMHC	0
NOX_GT	0
PT08_S3_NOX	0
NO2_GT	0
PT08_S4_NO2	0
PT08_S5_O3	0
T	0
RH	0

▼ Data Visualization

```
1 correlation_matrix = data.corr()
2 plt.figure(figsize=(12, 10))
3 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
4 plt.show()

<ipython-input-19-6c76563d8704>:1: FutureWarning: The default value of numeric_only
correlation_matrix = data.corr()
```



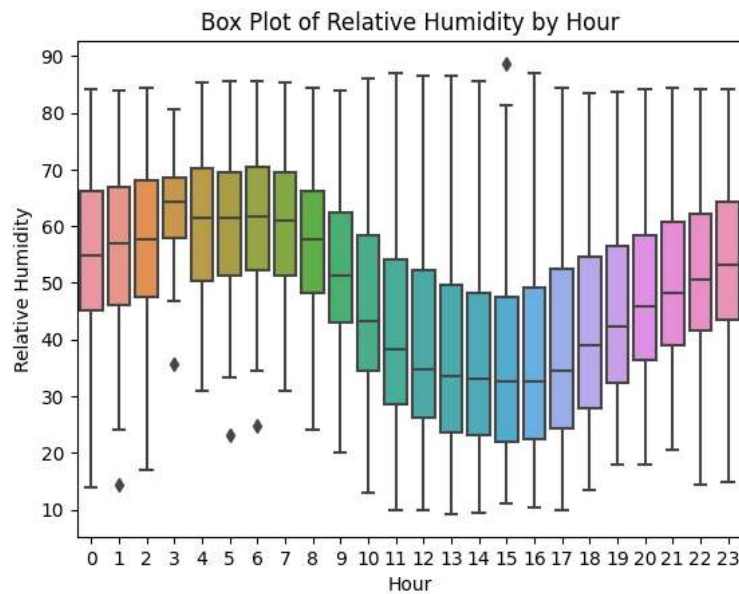
```
1 data['HOUR']=data['TIME'].apply(lambda x: int(x.split(':')[0]))
2 data.HOUR

0      18
1      19
2      20
3      21
4      22
..
9352   10
9353   11
9354   12
9355   13
9356   14
Name: HOUR, Length: 6941, dtype: int64
```

```

1 sns.boxplot(x='HOUR', y='RH', data=data)
2 plt.xlabel('Hour')
3 plt.ylabel('Relative Humidity')
4 plt.title('Box Plot of Relative Humidity by Hour')
5 plt.show()
6

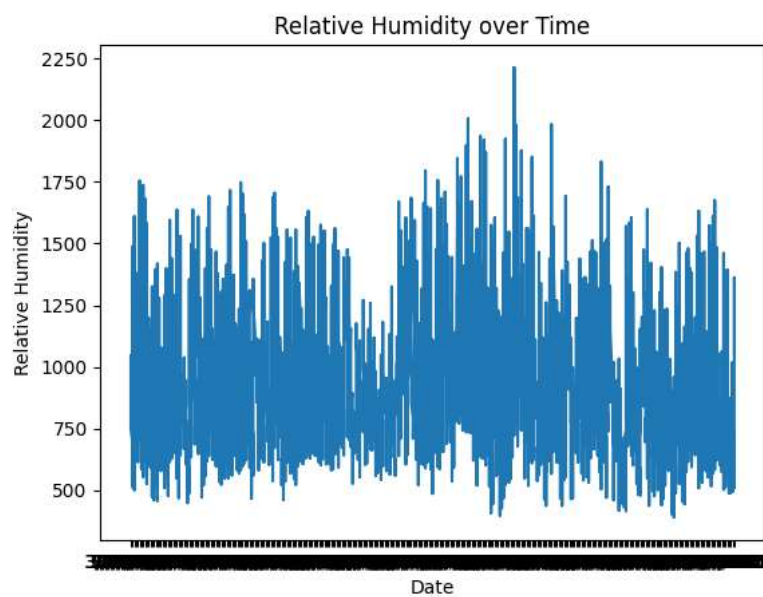
```



```

1 plt.plot(data['DATE'], data['PT08_S2_NMHC'])
2 plt.xlabel('Date')
3 plt.ylabel('Relative Humidity')
4 plt.title('Relative Humidity over Time')
5 plt.show()
6

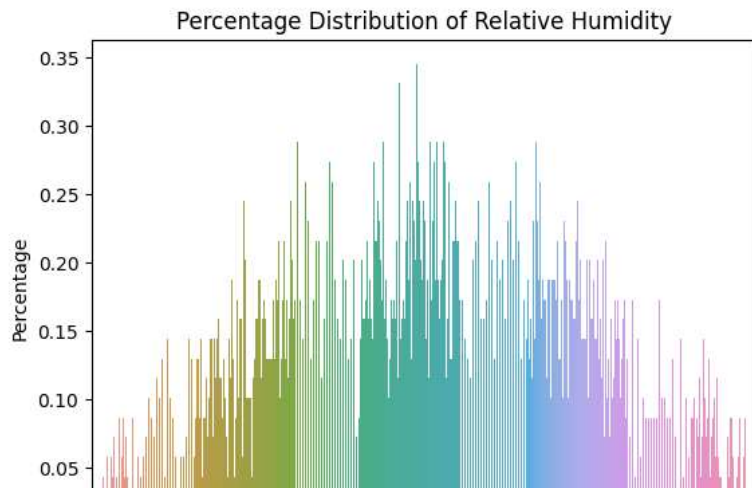
```



```

1 # Import the necessary libraries
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Plotting a bar plot of RH
6 sns.barplot(data=data, x='RH', y='RH', estimator=lambda x: len(x) / len(data) * 100)
7 plt.xlabel('Relative Humidity')
8 plt.ylabel('Percentage')
9 plt.title('Percentage Distribution of Relative Humidity')
10 plt.show()
11

```



## Linear Regression

```

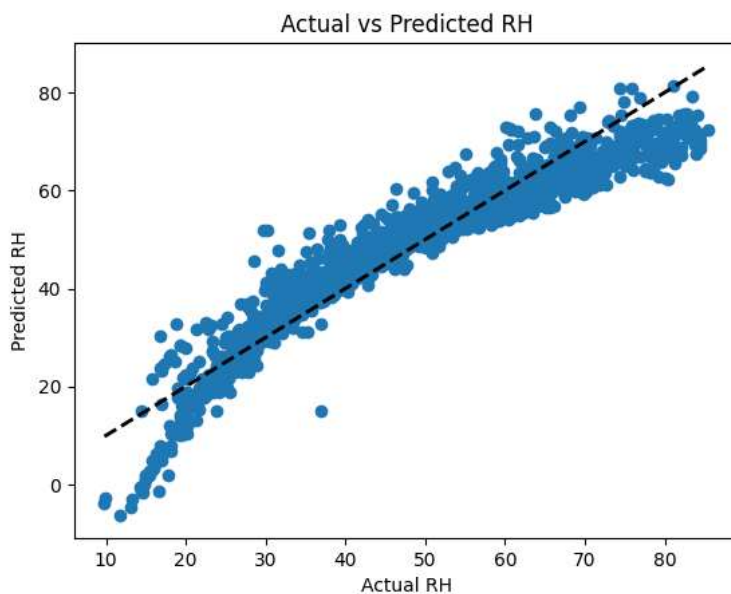
relative humidity

1 from sklearn.model_selection import train_test_split      #import train test split
2 from sklearn.linear_model import LinearRegression         #import linear regression package
3 from sklearn.metrics import mean_squared_error,mean_absolute_error  #import mean squared error and mean absolute error
4 from sklearn.metrics import r2_score

1 X = data.drop(['RH','DATE','TIME'], axis=1)
2 y = data['RH']
3
4 # Split the data into training and testing sets
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
6
7 # Create and train the linear regression model
8 linear_regression = LinearRegression()
9 linear_regression.fit(X_train, y_train)
10
11 # Make predictions on the test set
12 y_pred_lr = linear_regression.predict(X_test)

1 # Plotting actual vs predicted values
2 plt.scatter(y_test, y_pred_lr)
3 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
4 plt.xlabel('Actual RH')
5 plt.ylabel('Predicted RH')
6 plt.title('Actual vs Predicted RH')
7 plt.show()

```



```

1 mean_squared_error(y_test, y_pred_lr)

34.0802484967582

```

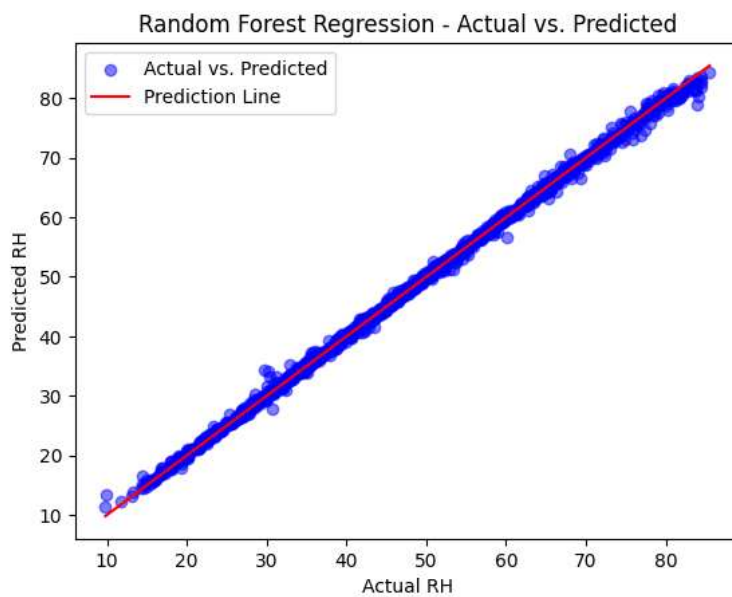
```
1 r2_score(y_test, y_pred_lr)

0.8899362175811059
```

## Random Forest Regression

```
1 from sklearn.ensemble import RandomForestRegressor
2 random_forest_model = RandomForestRegressor()
3 random_forest_model.fit(X_train, y_train)
4 y_pred_rf = random_forest_model.predict(X_test)

1 plt.scatter(y_test, y_pred_rf, color='blue', label='Actual vs. Predicted',alpha=0.5)
2 plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', label='Prediction Line')
3 plt.xlabel('Actual RH')
4 plt.ylabel('Predicted RH')
5 plt.title('Random Forest Regression - Actual vs. Predicted')
6 plt.legend()
7 plt.show()
```



```
1 r2_score(y_test, y_pred_rf)

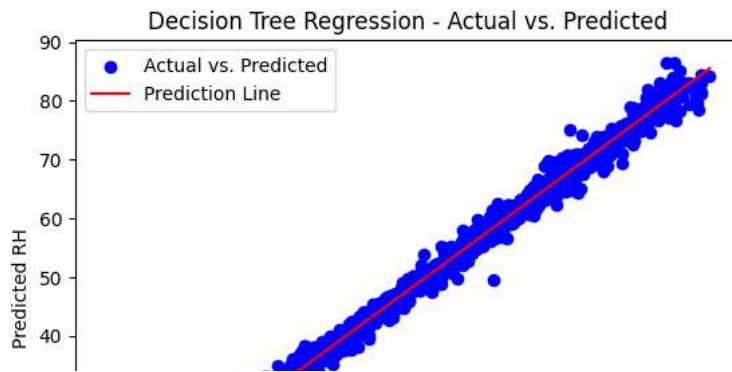
0.9985749187561072
```

## Decision Tree Regression

```
1 from sklearn.tree import DecisionTreeRegressor

1 decision_tree_model = DecisionTreeRegressor()
2 decision_tree_model.fit(X_train, y_train)
3 y_pred_dt = decision_tree_model.predict(X_test)

1 plt.scatter(y_test, y_pred_dt, color='blue', label='Actual vs. Predicted')
2 plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', label='Prediction Line')
3 plt.xlabel('Actual RH')
4 plt.ylabel('Predicted RH')
5 plt.title('Decision Tree Regression - Actual vs. Predicted')
6 plt.legend()
7 plt.show()
```



```
1 mean_squared_error(y_test, y_pred_rf)
```

```
0.4412634370050409
```

```
1 r2_score(y_test, y_pred_dt)
```

```
0.9940341491503423
```

## ▼ SVM

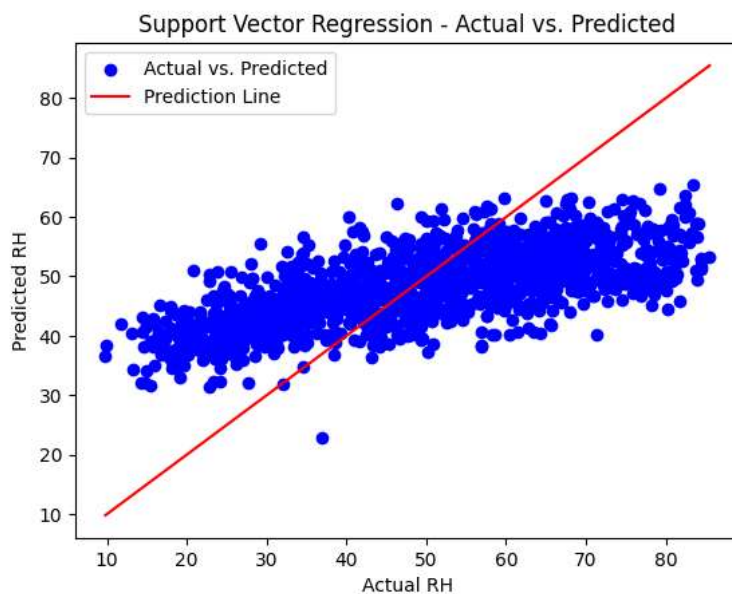
```
1 from sklearn.svm import SVR
```

```
1 svr_model = SVR()
2 svr_model.fit(X_train, y_train)
3 y_pred_svr = svr_model.predict(X_test)
```

```
1 r2_score(y_test, y_pred_svr)
```

```
0.3681904342526975
```

```
1 plt.scatter(y_test, y_pred_svr, color='blue', label='Actual vs. Predicted')
2 plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', label='Prediction Line')
3 plt.xlabel('Actual RH')
4 plt.ylabel('Predicted RH')
5 plt.title('Support Vector Regression - Actual vs. Predicted')
6 plt.legend()
7 plt.show()
```



## ▼ Lasso Regression

```
1 from sklearn.linear_model import Lasso
2 lasso_model = Lasso(alpha=0.1) # You can adjust the value of alpha
3 lasso_model.fit(X_train, y_train)
```

```

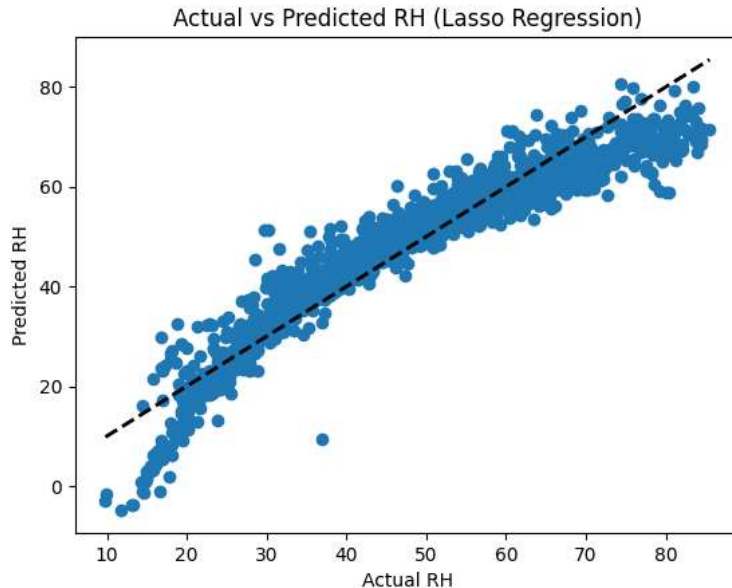
5 # Make predictions on the test set
6 y_pred_lm = lasso_model.predict(X_test)

1 r2_score(y_test, y_pred_lm)

0.8878969070776364

1 plt.scatter(y_test, y_pred_lm)
2 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
3 plt.xlabel('Actual RH')
4 plt.ylabel('Predicted RH')
5 plt.title('Actual vs Predicted RH (Lasso Regression)')
6 plt.show()
7

```



## Comparison

```

1 from sklearn.metrics import mean_squared_error, r2_score
2
3 # Calculate MSE for each models
4 mse_lr = mean_squared_error(y_test, y_pred_lr)
5 mse_rf = mean_squared_error(y_test, y_pred_rf)
6 mse_dt = mean_squared_error(y_test, y_pred_dt)
7 mse_svr = mean_squared_error(y_test, y_pred_svr)
8 mse_lm = mean_squared_error(y_test, y_pred_lm)
9
10 # Calculate R2 score for each model
11 r2_lr = r2_score(y_test, y_pred_lr)
12 r2_rf = r2_score(y_test, y_pred_rf)
13 r2_dt = r2_score(y_test, y_pred_dt)
14 r2_svr = r2_score(y_test, y_pred_svr)
15 r2_lm = r2_score(y_test, y_pred_lm)
16
17 # Create a dataframe to compare the metrics
18 metrics_df = pd.DataFrame({
19     'Model': ['Linear Regression', 'Random Forest', 'Decision Tree', 'SVR', 'Lasso Regression'],
20     'MSE': [mse_lr, mse_rf, mse_dt, mse_svr, mse_lm],
21     'R2 Score': [r2_lr, r2_rf, r2_dt, r2_svr, r2_lm]
22 })
23
24 # Display the dataframe
25 print(metrics_df)

```

	Model	MSE	R2 Score
0	Linear Regression	34.080248	0.889936
1	Random Forest	0.441263	0.998575
2	Decision Tree	1.847271	0.994034
3	SVR	195.634082	0.368190
4	Lasso Regression	34.711702	0.887897

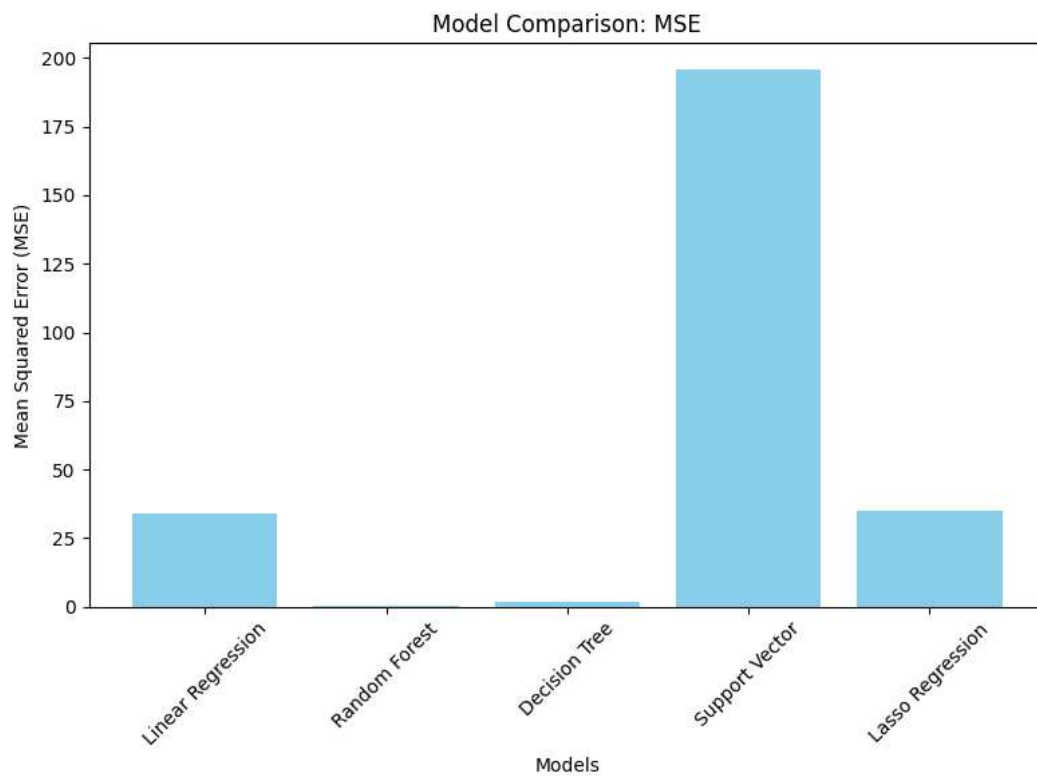
```

1 import numpy as np
2 import matplotlib.pyplot as plt

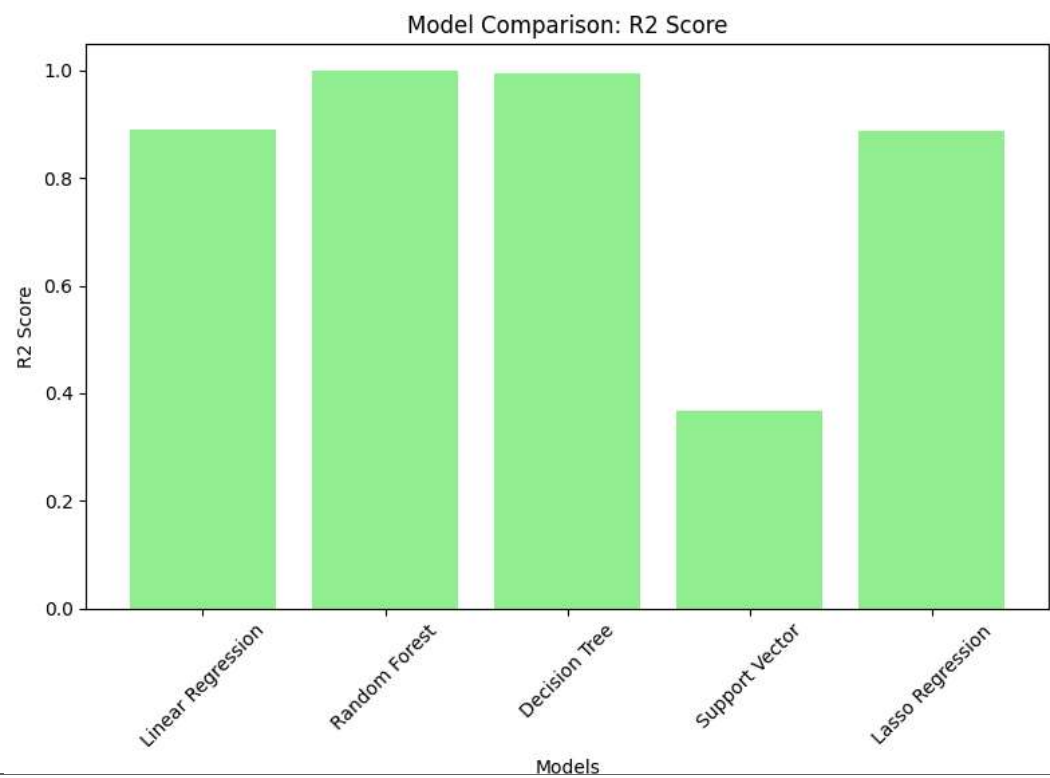
```



```
3
4 # Create lists of models, MSE values, and R2 scores
5 models = ['Linear Regression', 'Random Forest', 'Decision Tree', 'Support Vector', 'Lasso Regression']
6 mse_values = [mse_lr, mse_rf, mse_dt, mse_svr, mse_lm]
7 r2_scores = [r2_lr, r2_rf, r2_dt, r2_svr, r2_lm]
8
9 # Plot MSE comparison
10 plt.figure(figsize=(8, 6))
11 plt.bar(models, mse_values, color='skyblue')
12 plt.xlabel('Models')
13 plt.ylabel('Mean Squared Error (MSE)')
14 plt.title('Model Comparison: MSE')
15 plt.xticks(rotation=45)
16 plt.tight_layout()
17 plt.show()
18
19
```



```
1 # Plot R2 score comparison
2 plt.figure(figsize=(8, 6))
3 plt.bar(models, r2_scores, color='lightgreen')
4 plt.xlabel('Models')
5 plt.ylabel('R2 Score')
6 plt.title('Model Comparison: R2 Score')
7 plt.xticks(rotation=45)
8 plt.tight_layout()
9 plt.show()
10
```



✓ 0s completed at 10:09 PM

