# Data Collection and Preprocessing Phase

| Date | 9 July 2024 |
| --- | --- |
| Team ID | SWTID1720162737 |
| Project Title | Predicting Compressive Strength Of Concrete Using Machine Learning. |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
| --- | --- |
| Data Overview | **Dimensions:**<br><br>```[7]: data.shape```<br><br>```[7]: (1030, 9)```<br><br>**Descriptive statistics:**<br> |
| Univariate Analysis | ```#mean```<br>```np.mean(data)```<br><br>```269.444832793959``` |

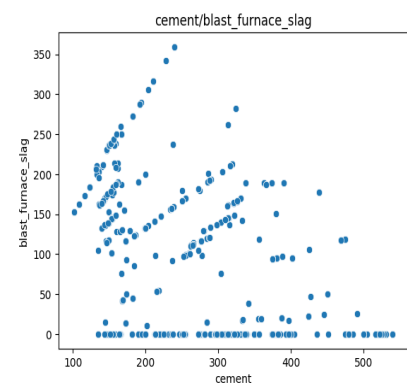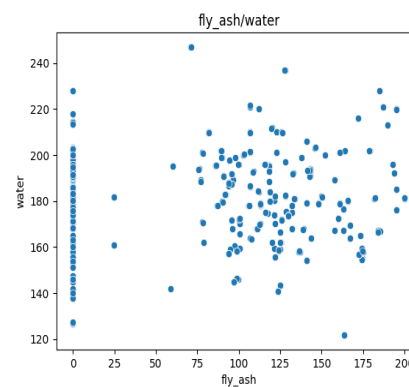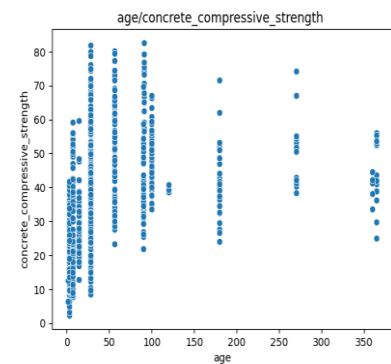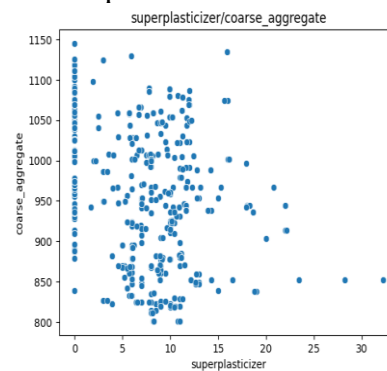| | |
|---|---|
| | ```
#median
np.median(data)

125.30000000000001

#mode
vals, counts = np.unique(data, return_counts=True)
max_count_index = np.argmax(counts)
mode_value = vals[max_count_index]
print("Mode:", mode_value)

Mode: 0.0
``` |
| Bivariate Analysis | **Correlation:**<br><br><br><br>**Scatter plots:**<br><br> |

| Multivariate Analysis |  |

| Outliers and Anomalies | With Outliers |
| --- | --- |
| |  |
| | Without outliers: |
| |  |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data | ```
data=pd.read_csv('concrete_data.csv')
data
``` <br><br> | | cement | blast_furnace_slag | fly_ash | water | superplasticizer | coarse_aggregate | fine_aggregate | age | concrete_compressive_strength |<br>|---|---|---|---|---|---|---|---|---|---|<br>| 0 | 540.0 | 0.0 | 0.0 | 162.0 | 2.5 | 1040.0 | 676.0 | 28 | 79.99 |<br>| 1 | 540.0 | 0.0 | 0.0 | 162.0 | 2.5 | 1055.0 | 676.0 | 28 | 61.89 |<br>| 2 | 332.5 | 142.5 | 0.0 | 228.0 | 0.0 | 932.0 | 594.0 | 270 | 40.27 |<br>| 3 | 332.5 | 142.5 | 0.0 | 228.0 | 0.0 | 932.0 | 594.0 | 365 | 41.05 |<br>| 4 | 198.6 | 132.4 | 0.0 | 192.0 | 0.0 | 978.4 | 825.5 | 360 | 44.30 | |
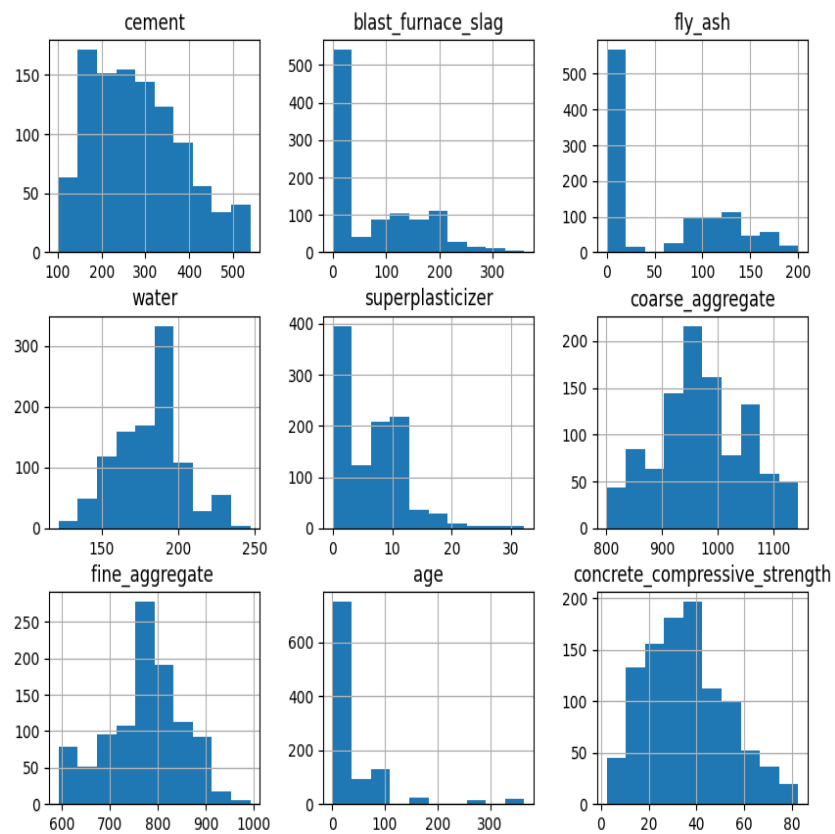| Handling Missing Data | ```
data['cement']=data['cement'].fillna(data['cement'].mode()[0])
data['blast_furnace_slag']=data['blast_furnace_slag'].fillna(data['blast_furnace_slag'].mode()[0])
data['fly_ash']=data['fly_ash'].fillna(data['fly_ash'].mode()[0])
data['water']=data['water'].fillna(data['water'].mode()[0])
data['coarse_aggregate']=data['coarse_aggregate'].fillna(data['coarse_aggregate'].mode()[0])
data['superplasticizer']=data['superplasticizer'].fillna(data['superplasticizer'].mode()[0])
data['age']=data['age'].fillna(data['age'].mode()[0])
data['concrete_compressive_strength']=data['concrete_compressive_strength'].fillna(data['concrete_compressive_strength'].mode()[0])
``` <br><br> But it is not necessary cause we don't have any missed values in the data set. |
| Data Transformation | ```
#Scaling on Independent variables
from sklearn.preprocessing import StandardScaler

scale=StandardScaler()

names=x.columns
names

Index(['cement', 'blast_furnace_slag', 'fly_ash', 'water', 'superplasticizer',
       'coarse_aggregate', 'fine_aggregate ', 'concrete_compressive_strength'],
      dtype='object')

scale.fit_transform(x)
``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | ```
data=filtered_data
data.shape

(1021, 9)
``` |