```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
sns.set_theme(color_codes=True)
pd.set_option('display.max_columns', None)


df = pd.read_csv('supply_chain_data.csv')
df.head()
```

| | Product type | SKU | Price | Availability | Number of products sold | Revenue generated | Customer demographics | Stock levels |
|---|---|---|---|---|---|---|---|---|
| 0 | haircare | SKU0 | 69.808006 | 55 | 802 | 8661.996792 | Non-binary | 58 |
| 1 | skincare | SKU1 | 14.843523 | 95 | 736 | 7460.900065 | Female | 53 |
| 2 | haircare | SKU2 | 11.319683 | 34 | 8 | 9577.749626 | Unknown | 1 |
| 3 | skincare | SKU3 | 61.163343 | 68 | 83 | 7766.836426 | Non-binary | 23 |
| 4 | skincare | SKU4 | 4.805496 | 26 | 871 | 2686.505152 | Non-binary | 5 |

## Data Preprocessing Part 1

```
df.dtypes
```

```
Product type             object
SKU                      object
Price                    float64
Availability             int64
Number of products sold  int64
Revenue generated        float64
Customer demographics    object
Stock levels             int64
Lead times               int64
Order quantities         int64
Shipping times           int64
Shipping carriers        object
Shipping costs           float64
Supplier name            object
Location                 object
Lead time                int64
Production volumes       int64
Manufacturing lead time  int64
Manufacturing costs      float64
```

```
Inspection results          object
Defect rates                float64
Transportation modes        object
Routes                      object
Costs                       float64
dtype: object
```

```
df.select_dtypes(include='object').nunique()
```

```
Product type            3
SKU                   100
Customer demographics   4
Shipping carriers       3
Supplier name           5
Location                5
Inspection results      3
Transportation modes    4
Routes                  3
dtype: int64
```

```
df.shape
```

```
(100, 24)
```

```
#Drop SKU Column because this is just supply chain id
df.drop(columns=['SKU'], inplace=True)
df.shape
```

```
(100, 23)
```

```
df.nunique()
```

```
Product type              3
Price                   100
Availability             63
Number of products sold  96
Revenue generated       100
Customer demographics     4
Stock levels             65
Lead times               29
Order quantities         61
Shipping times           10
Shipping carriers         3
Shipping costs          100
Supplier name             5
Location                  5
Lead time                29
Production volumes       96
Manufacturing lead time  30
Manufacturing costs     100
Inspection results        3
Defect rates            100
Transportation modes      4
Routes                    3
Costs                   100
dtype: int64
```
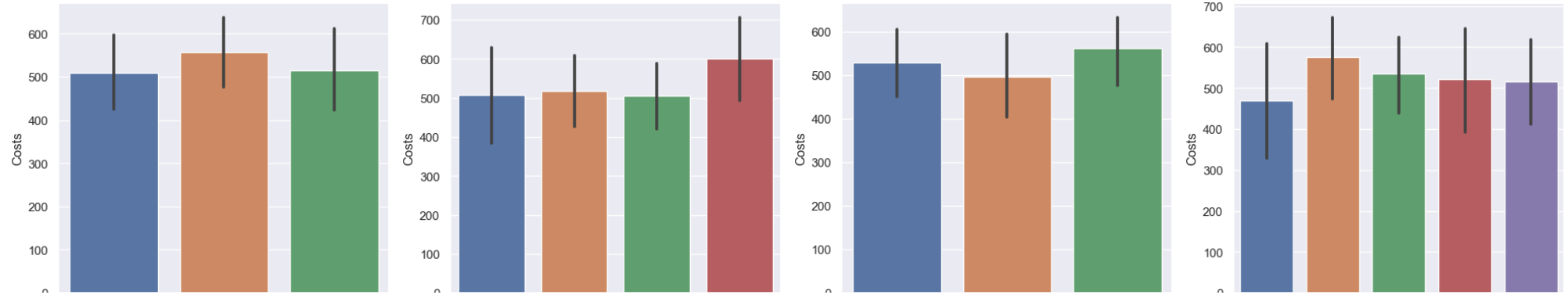
# ▾ Exploratory Data Analysis

```python
# list of categorical variables to plot
cat_vars = ['Product type', 'Customer demographics', 'Shipping carriers', 'Supplier name', 'Location',
            'Inspection results', 'Transportation modes', 'Routes']

# create figure with subplots
fig, axs = plt.subplots(nrows=2, ncols=4, figsize=(20, 10))
axs = axs.flatten()

# create barplot for each categorical variable
for i, var in enumerate(cat_vars):
    sns.barplot(x=var, y='Costs', data=df, ax=axs[i], estimator=np.mean)
    axs[i].set_xticklabels(axs[i].get_xticklabels(), rotation=90)

# adjust spacing between subplots
fig.tight_layout()

# show plot
plt.show()
```

```
num_vars = ['Price', 'Availability', 'Number of products sold', 'Stock levels',
            'Lead times', 'Order quantities', 'Shipping times', 'Shipping costs',
            'Lead time', 'Production volumes', 'Manufacturing lead time', 'Manufacturing costs',
            'Defect rates']

fig, axs = plt.subplots(nrows=3, ncols=5, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.boxplot(x=var, data=df, ax=axs[i])

# remove the 14th subplot
fig.delaxes(axs[13])
# remove the 15th subplot
fig.delaxes(axs[14])

fig.tight_layout()
plt.show()
```
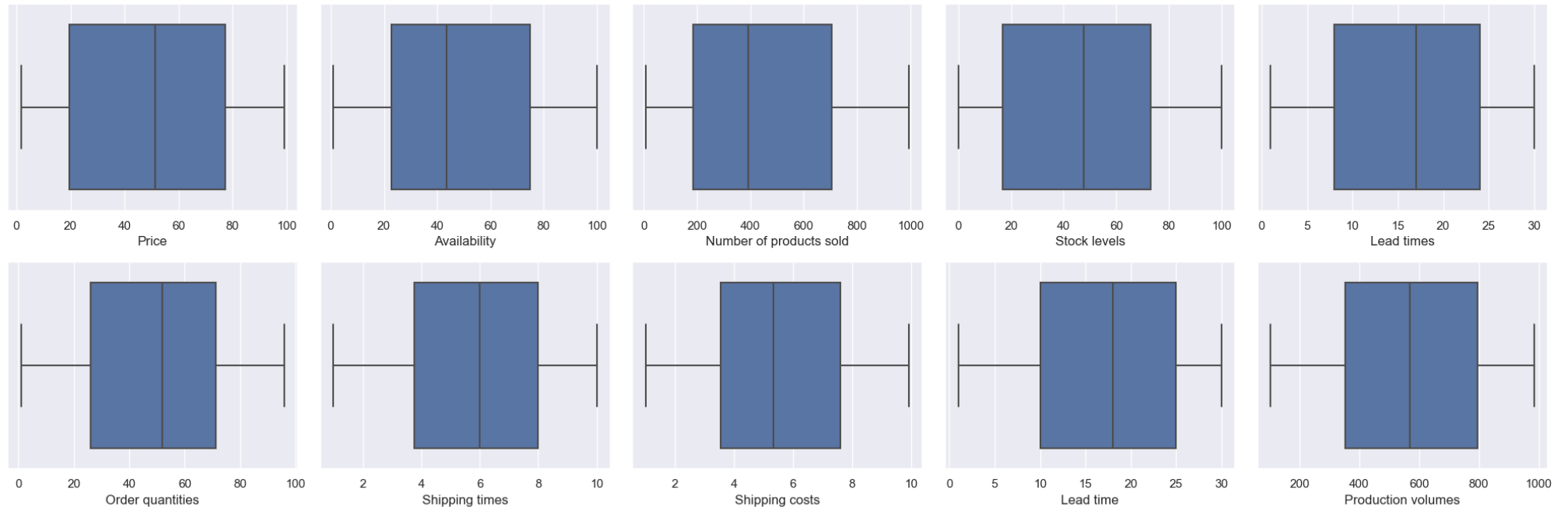
```
num_vars = ['Price', 'Availability', 'Number of products sold', 'Stock levels',
            'Lead times', 'Order quantities', 'Shipping times', 'Shipping costs',
            'Lead time', 'Production volumes', 'Manufacturing lead time', 'Manufacturing costs',
            'Defect rates']

fig, axs = plt.subplots(nrows=3, ncols=5, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.violinplot(x=var, data=df, ax=axs[i])

# remove the 14th subplot
fig.delaxes(axs[13])
# remove the 15th subplot
fig.delaxes(axs[14])

fig.tight_layout()
plt.show()
```
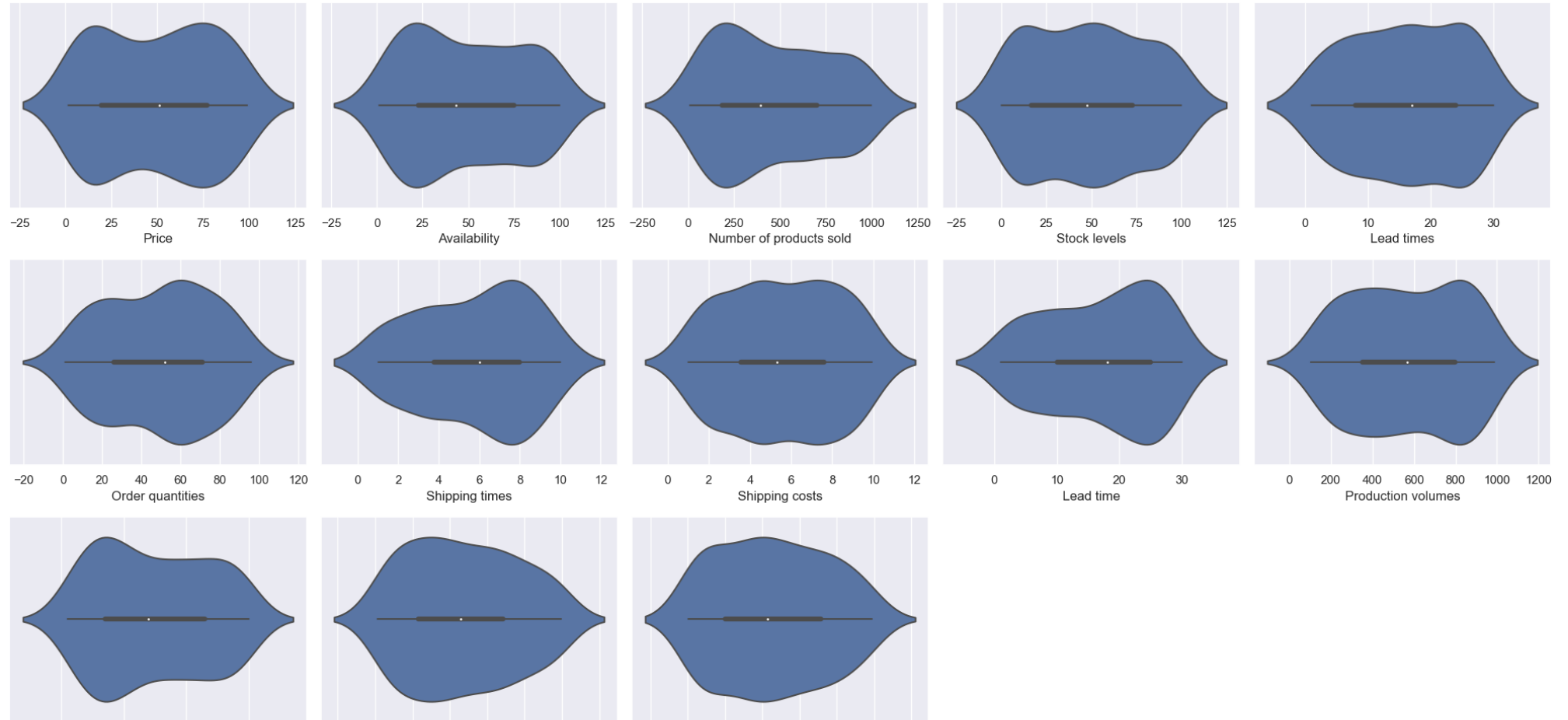
```
num_vars = ['Price', 'Availability', 'Number of products sold', 'Stock levels',
            'Lead times', 'Order quantities', 'Shipping times', 'Shipping costs',
            'Lead time', 'Production volumes', 'Manufacturing lead time', 'Manufacturing costs',
            'Defect rates']

fig, axs = plt.subplots(nrows=3, ncols=5, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.scatterplot(x=var, y='Costs', hue='Routes', data=df, ax=axs[i])

# remove the 14th subplot
fig.delaxes(axs[13])
# remove the 15th subplot
fig.delaxes(axs[14])

fig.tight_layout()
```

```
plt.show()
```



```
num_vars = ['Price', 'Availability', 'Number of products sold', 'Stock levels',
            'Lead times', 'Order quantities', 'Shipping times', 'Shipping costs',
            'Lead time', 'Production volumes', 'Manufacturing lead time', 'Manufacturing costs',
            'Defect rates']
```
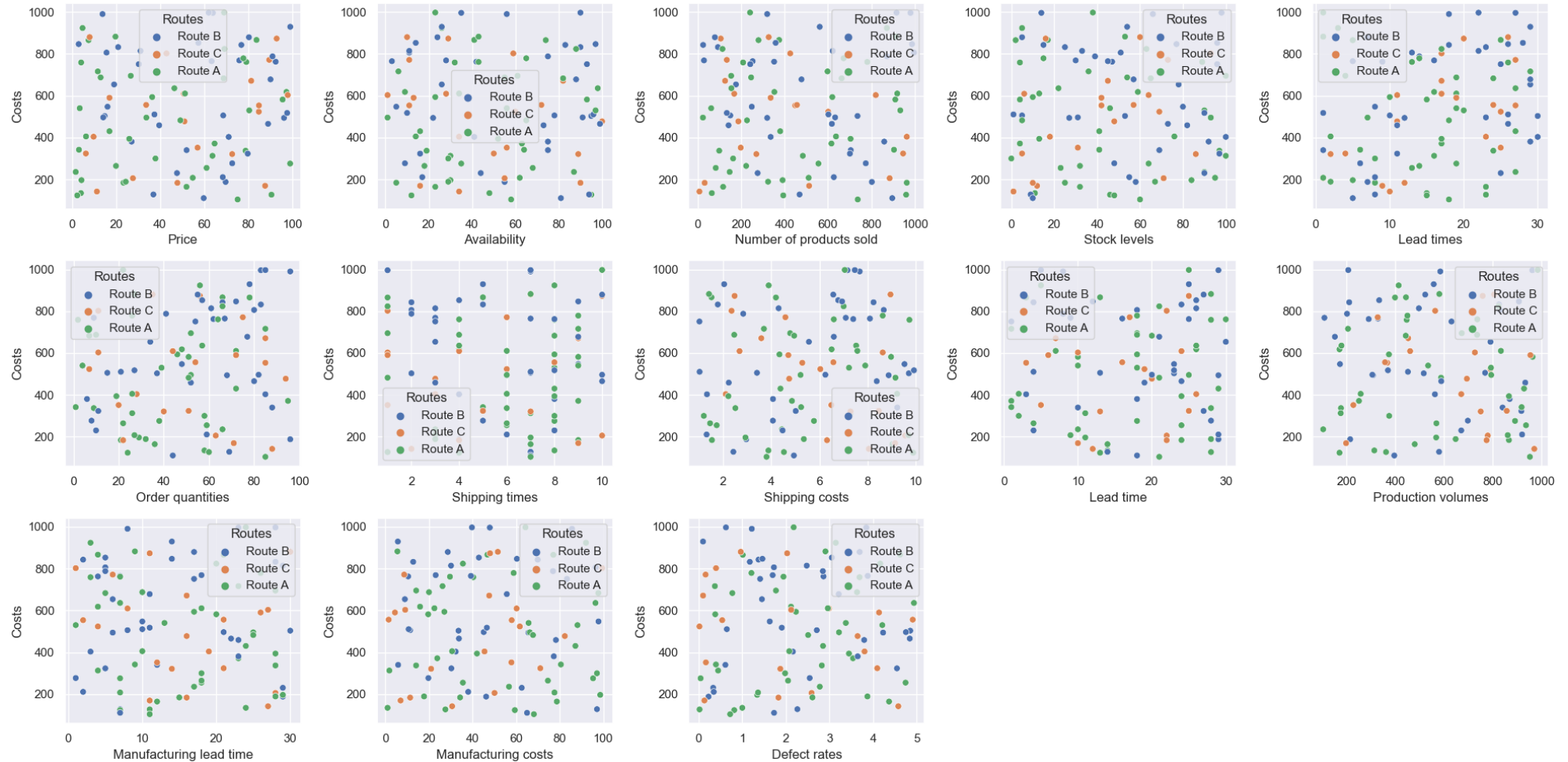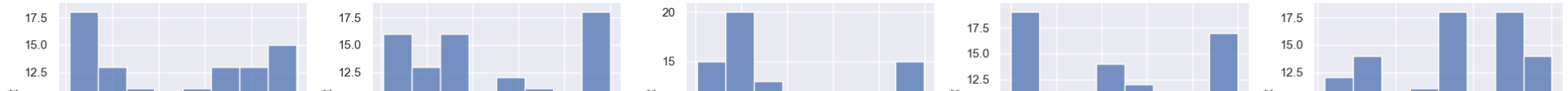
```
fig, axs = plt.subplots(nrows=3, ncols=5, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.histplot(x=var, data=df, ax=axs[i])

# remove the 14th subplot
fig.delaxes(axs[13])
# remove the 15th subplot
fig.delaxes(axs[14])

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

## Data Preprocessing Part 2

```
#Check the missing value
check_missing = df.isnull().sum() * 100 / df.shape[0]
check_missing[check_missing > 0].sort_values(ascending=False)
```

```
Series([], dtype: float64)
```

## Label Encoding for Object datatypes

```
# Loop over each column in the DataFrame where dtype is 'object'
for col in df.select_dtypes(include=['object']).columns:

    # Print the column name and the unique values
    print(f"{col}: {df[col].unique()}")
```

```
Product type: ['haircare' 'skincare' 'cosmetics']
Customer demographics: ['Non-binary' 'Female' 'Unknown' 'Male']
Shipping carriers: ['Carrier B' 'Carrier A' 'Carrier C']
Supplier name: ['Supplier 3' 'Supplier 1' 'Supplier 5' 'Supplier 4' 'Supplier 2']
Location: ['Mumbai' 'Kolkata' 'Delhi' 'Bangalore' 'Chennai']
Inspection results: ['Pending' 'Fail' 'Pass']
Transportation modes: ['Road' 'Air' 'Rail' 'Sea']
Routes: ['Route B' 'Route C' 'Route A']
```

```
from sklearn import preprocessing

# Loop over each column in the DataFrame where dtype is 'object'
for col in df.select_dtypes(include=['object']).columns:

    # Initialize a LabelEncoder object
    label_encoder = preprocessing.LabelEncoder()

    # Fit the encoder to the unique values in the column
    label_encoder.fit(df[col].unique())

    # Transform the column using the encoder
    df[col] = label_encoder.transform(df[col])

    # Print the column name and the unique encoded values
    print(f"{col}: {df[col].unique()}")
```

```
Product type: [1 2 0]
Customer demographics: [2 0 3 1]
Shipping carriers: [1 0 2]
```

```
Supplier name: [2 0 4 3 1]
Location: [4 3 2 0 1]
Inspection results: [2 0 1]
Transportation modes: [2 0 1 3]
Routes: [1 2 0]
```

```
df.dtypes
```

```
Product type                int32
Price                     float64
Availability                int64
Number of products sold     int64
Revenue generated         float64
Customer demographics       int32
Stock levels                int64
Lead times                  int64
Order quantities            int64
Shipping times              int64
Shipping carriers           int32
Shipping costs            float64
Supplier name               int32
Location                    int32
Lead time                   int64
Production volumes          int64
Manufacturing lead time     int64
Manufacturing costs       float64
Inspection results          int32
Defect rates              float64
Transportation modes        int32
Routes                      int32
Costs                     float64
dtype: object
```

# There's no outlier so we dont have to remove it

## ▾ Correlation Heatmap

```
#Correlation Heatmap
plt.figure(figsize=(20, 16))
sns.heatmap(df.corr(), fmt='.2g', annot=True)
```

```
<AxesSubplot:>
```

| | Product type | Price | Availability | Number of products sold | Revenue generated | Customer demographics | Stock levels | Lead times | Order quantities | Shipping times | Shipping carriers | Shipping costs | Supplier name | Location | Lead time | Production volumes | Manufacturing lead time | Manufacturing costs | Inspection results | Defect rates | Transportation modes | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Product type | 1 | -0.12 | 0.011 | 0.1 | -0.0035 | -0.015 | -0.23 | 0.064 | 0.031 | -0.18 | -0.1 | -0.18 | -0.093 | -0.042 | 0.18 | 0.19 | -0.0025 | 0.077 | 0.033 | 0.1 | -0.074 | 0.0036 | 0.071 |
| Price | -0.12 | 1 | 0.019 | 0.0057 | 0.038 | 0.14 | 0.078 | 0.045 | 0.096 | 0.072 | 0.2 | 0.059 | -0.15 | -0.046 | 0.15 | -0.12 | -0.3 | -0.18 | -0.061 | -0.15 | 0.009 | 0.15 | 0.089 |
| Availability | 0.011 | 0.019 | 1 | 0.087 | -0.075 | -0.031 | -0.026 | 0.17 | 0.14 | -0.051 | -0.0015 | -0.044 | 0.13 | -0.18 | -0.16 | 0.05 | 0.065 | 0.13 | -0.045 | 0.041 | 0.03 | 0.018 | -0.027 |
| Number of products sold | 0.1 | 0.0057 | 0.087 | 1 | -0.0016 | 0.015 | 0.022 | -0.046 | 0.016 | 0.087 | -0.24 | 0.044 | 0.033 | 0.14 | 0.041 | 0.19 | -0.049 | 0.034 | 0.13 | -0.083 | 0.076 | -0.053 | -0.037 |
| Revenue generated | -0.0035 | 0.038 | -0.075 | -0.0016 | 1 | -0.14 | -0.16 | -0.057 | 0.029 | -0.11 | 0.18 | -0.073 | -0.016 | 0.034 | -0.014 | -0.037 | 0.014 | -0.21 | 0.15 | -0.13 | -0.053 | -0.0021 | 0.027 |
| Customer demographics | -0.015 | 0.14 | -0.031 | 0.015 | -0.14 | 1 | 0.052 | -0.062 | 0.12 | 0.0095 | 0.016 | 0.037 | 0.048 | 0.13 | 0.015 | -0.0074 | 0.078 | -0.025 | 0.11 | 0.049 | -0.043 | 0.088 | -0.056 |
| Stock levels | -0.23 | 0.078 | -0.026 | 0.022 | -0.16 | 0.052 | 1 | 0.073 | -0.11 | -0.095 | -0.098 | 0.073 | 0.13 | 0.038 | 0.068 | 0.044 | -0.051 | 0.033 | 0.076 | -0.15 | -0.049 | 0.0066 | -0.012 |
| Lead times | 0.064 | 0.045 | 0.17 | -0.046 | -0.057 | -0.062 | 0.073 | 1 | 0.11 | -0.045 | 0.088 | -0.12 | -0.059 | -0.061 | -0.0028 | -0.15 | 0.0034 | -0.024 | -0.12 | 0.016 | -0.17 | 0.093 | 0.24 |
| Order quantities | 0.031 | 0.096 | 0.14 | 0.016 | 0.029 | 0.12 | -0.11 | 0.11 | 1 | -0.0026 | 0.096 | 0.0043 | -0.02 | 0.028 | -0.086 | -0.087 | 0.11 | -0.027 | -0.084 | 0.019 | -0.025 | 0.15 | 0.17 |
| Shipping times | -0.18 | 0.072 | -0.051 | 0.087 | -0.11 | 0.0095 | -0.095 | -0.045 | -0.0026 | 1 | -0.013 | 0.045 | 0.0038 | 0.021 | -0.022 | -0.06 | -0.017 | 0.029 | -0.092 | -0.037 | 0.12 | -0.11 | -0.046 |
| Shipping carriers | -0.1 | 0.2 | -0.0015 | -0.24 | 0.18 | 0.016 | -0.098 | 0.088 | 0.096 | -0.013 | 1 | 0.0065 | 0.038 | 0.17 | -0.05 | -0.13 | -0.064 | -0.096 | -0.12 | -0.034 | 0.021 | 0.039 | 0.094 |
| Shipping costs | -0.18 | 0.059 | -0.044 | 0.044 | -0.073 | 0.037 | 0.073 | -0.12 | 0.0043 | 0.045 | 0.0065 | 1 | 0.026 | 0.12 | 0.03 | -0.098 | -0.0057 | 0.006 | 0.12 | 0.083 | -0.12 | 0.072 | 0.052 |
| Supplier name | -0.093 | -0.15 | 0.13 | 0.033 | -0.016 | 0.048 | 0.13 | -0.059 | -0.02 | 0.0038 | 0.038 | 0.026 | 1 | 0.0041 | 0.073 | 0.047 | 0.13 | 0.1 | -0.061 | 0.18 | 0.17 | -0.14 | -0.054 |
| Location | -0.042 | -0.046 | -0.18 | 0.14 | 0.034 | 0.13 | 0.038 | -0.061 | 0.028 | 0.021 | 0.17 | 0.12 | 0.0041 | 1 | -0.017 | 0.18 | 0.3 | -0.29 | 0.13 | -0.035 | 0.0013 | 0.055 | -0.25 |
| Lead time | 0.18 | 0.15 | -0.16 | 0.041 | -0.014 | 0.015 | 0.068 | -0.0028 | -0.086 | -0.022 | -0.05 | 0.03 | 0.073 | -0.017 | 1 | 0.21 | 0.027 | -0.12 | 0.092 | 0.3 | 0.037 | 0.028 | 0.045 |
| Production volumes | 0.19 | -0.12 | 0.05 | 0.19 | -0.037 | -0.0074 | 0.044 | -0.15 | -0.087 | -0.06 | -0.13 | -0.098 | 0.047 | 0.18 | 0.21 | 1 | 0.18 | 0.052 | 0.16 | 0.12 | 0.0077 | 0.066 | -0.075 |
| Manufacturing lead time | -0.0025 | -0.3 | 0.065 | -0.049 | 0.014 | 0.078 | -0.051 | 0.0034 | 0.11 | -0.017 | -0.064 | -0.0057 | 0.13 | 0.3 | 0.027 | 0.18 | 1 | -0.16 | 0.27 | 0.14 | -0.095 | -0.0079 | -0.074 |
| Manufacturing costs | 0.077 | -0.18 | 0.13 | 0.034 | -0.21 | -0.025 | 0.033 | -0.024 | -0.027 | 0.029 | -0.096 | 0.006 | 0.1 | -0.29 | -0.12 | 0.052 | -0.16 | 1 | -0.13 | -0.0078 | 0.023 | -0.12 | -0.014 |
| Inspection results | 0.033 | -0.061 | -0.045 | 0.13 | 0.15 | 0.11 | 0.076 | -0.12 | -0.084 | -0.092 | -0.12 | 0.12 | -0.061 | 0.13 | 0.092 | 0.16 | 0.27 | -0.13 | 1 | -0.12 | -0.031 | 0.017 | 0.013 |
| Defect rates | 0.1 | -0.15 | 0.041 | -0.083 | -0.13 | 0.049 | -0.15 | 0.016 | 0.019 | -0.037 | -0.034 | 0.083 | 0.18 | -0.035 | 0.3 | 0.12 | 0.14 | -0.0078 | -0.12 | 1 | 0.15 | -0.065 | 0.032 |
| Transportation modes | -0.074 | 0.009 | 0.03 | 0.076 | -0.053 | -0.043 | -0.049 | -0.17 | -0.025 | 0.12 | 0.021 | -0.12 | 0.17 | 0.0013 | 0.037 | 0.0077 | -0.095 | 0.023 | -0.031 | 0.15 | 1 | -0.0062 | -0.15 |

## Train test Split

```python
X = df.drop('Costs', axis=1)
y = df['Costs']
```

```python
#test size 20% and train size 80%
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2,random_state=0)
```

## Decision Tree Regressor

```python
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import load_boston


# Create a DecisionTreeRegressor object
dtree = DecisionTreeRegressor()

# Define the hyperparameters to tune and their values
param_grid = {
    'max_depth': [2, 4, 6, 8],
    'min_samples_split': [2, 4, 6, 8],
    'min_samples_leaf': [1, 2, 3, 4],
    'max_features': ['auto', 'sqrt', 'log2'],
    'random_state': [0, 7, 42]
}

# Create a GridSearchCV object
grid_search = GridSearchCV(dtree, param_grid, cv=5, scoring='neg_mean_squared_error')

# Fit the GridSearchCV object to the data
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print(grid_search.best_params_)
```

```
    {'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 3, 'min_samples_split': 2, 'random_state': 0}
```

```python
from sklearn.tree import DecisionTreeRegressor
dtree = DecisionTreeRegressor(random_state=0, max_depth=2, max_features='sqrt', min_samples_leaf=3, min_samples_split=2)
dtree.fit(X_train, y_train)
```

```
    DecisionTreeRegressor(max_depth=2, max_features='sqrt', min_samples_leaf=3,
                          random_state=0)
```

```python
from sklearn import metrics
from sklearn.metrics import mean_absolute_percentage_error
import math
y_pred = dtree.predict(X_test)
mae = metrics.mean_absolute_error(y_test, y_pred)
mape = mean_absolute_percentage_error(y_test, y_pred)
mse = metrics.mean_squared_error(y_test, y_pred)
r2 = metrics.r2_score(y_test, y_pred)
rmse = math.sqrt(mse)

print('MAE is {}'.format(mae))
print('MAPE is {}'.format(mape))
print('MSE is {}'.format(mse))
print('R2 score is {}'.format(r2))
print('RMSE score is {}'.format(rmse))
```

```
    MAE is 248.4413893861546
    MAPE is 0.5893818876444419
    MSE is 72806.47766651674
    R2 score is -0.08647889188367719
    RMSE score is 269.8267549123266
```

```python
imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": dtree.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Feature Importance Each Attributes (Decision Tree Regressor)', fontsize=18)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```

Feature Importance Each Attributes (Decision Tree Regressor)

```
import shap
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

## Random Forest Regressor

```python
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import GridSearchCV

# Create a Random Forest Regressor object
rf = RandomForestRegressor()

# Define the hyperparameter grid
param_grid = {
    'max_depth': [3, 5, 7, 9],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'max_features': ['auto', 'sqrt'],
    'random_state': [0, 7, 42]
}

# Create a GridSearchCV object
grid_search = GridSearchCV(rf, param_grid, cv=5, scoring='r2')

# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print("Best hyperparameters: ", grid_search.best_params_)
```

```
    Best hyperparameters:  {'max_depth': 3, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_split': 5, 'random_state': 0}
```

```python
from sklearn.ensemble import RandomForestRegressor
rf = RandomForestRegressor(random_state=0, max_depth=3, min_samples_split=5, min_samples_leaf=2,
                           max_features='sqrt')
rf.fit(X_train, y_train)
```

```
    RandomForestRegressor(max_depth=3, max_features='sqrt', min_samples_leaf=2,
                          min_samples_split=5, random_state=0)
```

```python
from sklearn import metrics
from sklearn.metrics import mean_absolute_percentage_error
import math
y_pred = rf.predict(X_test)
mae = metrics.mean_absolute_error(y_test, y_pred)
mape = mean_absolute_percentage_error(y_test, y_pred)
mse = metrics.mean_squared_error(y_test, y_pred)
r2 = metrics.r2_score(y_test, y_pred)
rmse = math.sqrt(mse)

print('MAE is {}'.format(mae))
print('MAPE is {}'.format(mape))
print('MSE is {}'.format(mse))
print('R2 score is {}'.format(r2))
print('RMSE score is {}'.format(rmse))
```

```
MAE is 247.33969719962744
MAPE is 0.6029768224226728
MSE is 71899.28833186119
R2 score is -0.07294105713825938
RMSE score is 268.14042651540103
```

```python
imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": dtree.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Feature Importance Each Attributes (Random Forest Regressor)', fontsize=18)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```
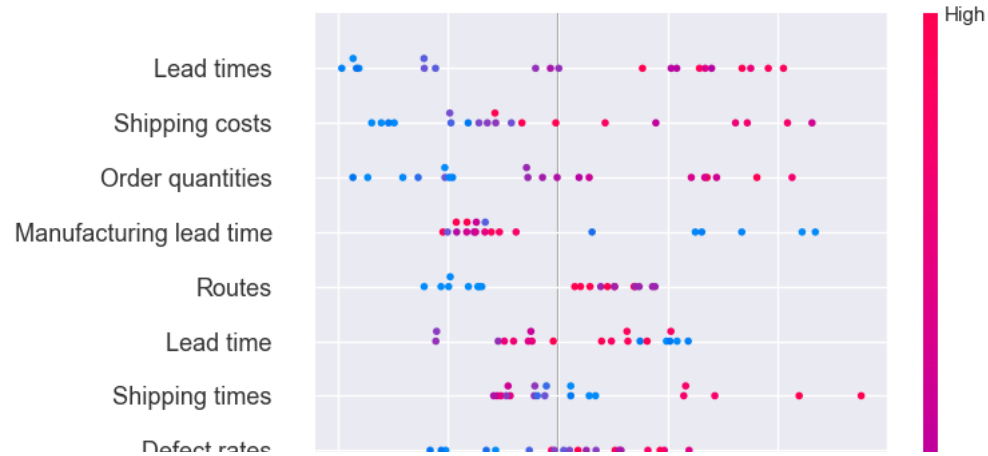
Feature Importance Each Attributes (Random Forest Regressor)

```
import shap
explainer = shap.TreeExplainer(rf)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

## AdaBoost Regressor



```
from sklearn.ensemble import AdaBoostRegressor
from sklearn.model_selection import GridSearchCV

# Create an AdaBoost Regressor object
ada = AdaBoostRegressor()

# Define the hyperparameter grid
param_grid = {
    'n_estimators': [50, 100, 150, 200],
    'learning_rate': [0.01, 0.1, 1],
    'loss': ['linear', 'square', 'exponential'],
    'random_state': [0, 7, 42]
}

# Create a GridSearchCV object
grid_search = GridSearchCV(ada, param_grid, cv=5, scoring='r2')

# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print("Best hyperparameters: ", grid_search.best_params_)
```

```
    Best hyperparameters:  {'learning_rate': 1, 'loss': 'linear', 'n_estimators': 50, 'random_state': 7}
```

```
from sklearn.ensemble import AdaBoostRegressor
ada = AdaBoostRegressor(random_state=7, n_estimators=50, learning_rate=1, loss='linear')
ada.fit(X_train, y_train)
```

```
    AdaBoostRegressor(learning_rate=1, random_state=7)
```
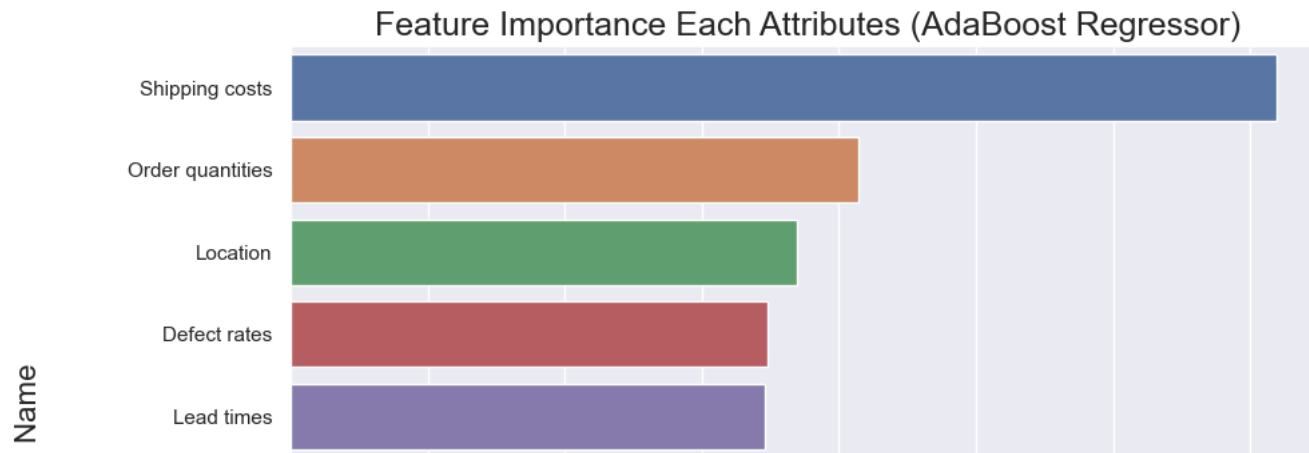
```python
from sklearn import metrics
from sklearn.metrics import mean_absolute_percentage_error
import math
y_pred = ada.predict(X_test)
mae = metrics.mean_absolute_error(y_test, y_pred)
mape = mean_absolute_percentage_error(y_test, y_pred)
mse = metrics.mean_squared_error(y_test, y_pred)
r2 = metrics.r2_score(y_test, y_pred)
rmse = math.sqrt(mse)

print('MAE is {}'.format(mae))
print('MAPE is {}'.format(mape))
print('MSE is {}'.format(mse))
print('R2 score is {}'.format(r2))
print('RMSE score is {}'.format(rmse))
```
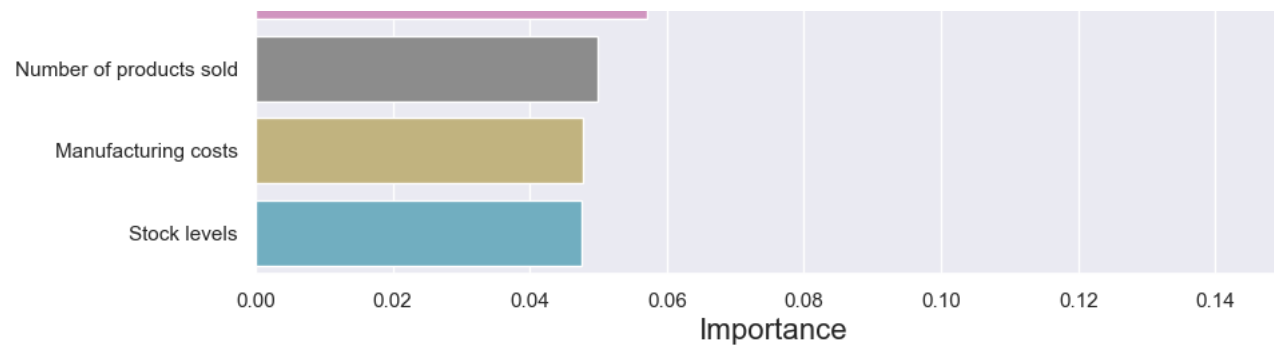
```
    MAE is 255.28612180541396
    MAPE is 0.5859844238678936
    MSE is 78800.74415400976
    R2 score is -0.17593032834538103
    RMSE score is 280.71470241868303
```

```python
imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": ada.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Feature Importance Each Attributes (AdaBoost Regressor)', fontsize=18)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```

## Feature Importance Each Attributes (AdaBoost Regressor)



All of the Algorithms got bad R2 Score and MAPE Score even with hyperparameter tuning because we only have 100 data and the distribution is spread



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.