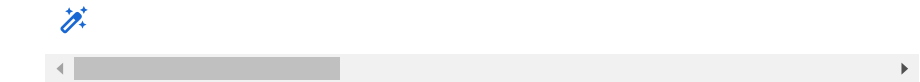


```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
sns.set_theme(color_codes=True)
pd.set_option("display.max_columns",None)
```

```
df=pd.read_csv("/content/supply_chain_data.csv")
df
```

	Product type	SKU	Price	Availability	Number of products sold	Revenue generated	Customer demographics	Stock level
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
...
95	haircare	SKU95	77.903927	65	672	7386.363944	Unknown	58
96	cosmetics	SKU96	24.423131	29	324	7698.424766	Non-binary	53
97	haircare	SKU97	3.526111	56	62	4370.916580	Male	1
98	skincare	SKU98	19.754605	43	913	8525.952560	Female	30
99	haircare	SKU99	68.517833	17	627	9185.185829	Unknown	10

100 rows × 24 columns



```
df.head()
```

	Product type	SKU	Price	Availability	Number of products sold	Revenue generated	Customer demographics	Stock levels	Lead times	Order quantities
0	haircare	SKU0	69.808006	55	802	8661.996792	Non-binary	58	7	5
1	skincare	SKU1	14.843523	95	736	7460.900065	Female	53	30	4
2	haircare	SKU2	11.319683	34	8	9577.749626	Unknown	1	10	8
3	skincare	SKU3	61.163343	68	83	7766.836426	Non-binary	23	13	4
4	skincare	SKU4	4.805496	26	871	2686.505152	Non-binary	5	3	4



Data preprocessing Part 1

```
df.dtypes
df.select_dtypes(include="object").nunique()

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

```
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

```
#list of categorical variables to plot
```

```
cat_vars=["Product type","Customer demographics","Shipping carriers","Supplier name","Location","Inspection results","Transportation mode"]
```

```
#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)
```

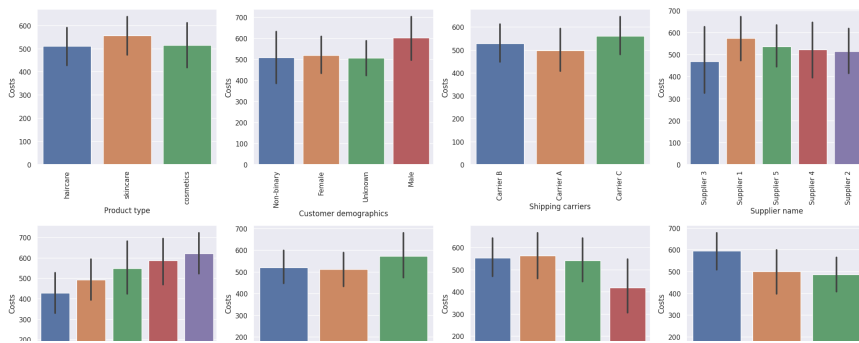
```
    axs[i].set_xticklabels(axs[i].get_xticklabels(),rotation=90)
```

```
#adjust spacing between subplots
```

```
fig.tight_layout()
```

```
#show plot
```

```
plt.show()
```

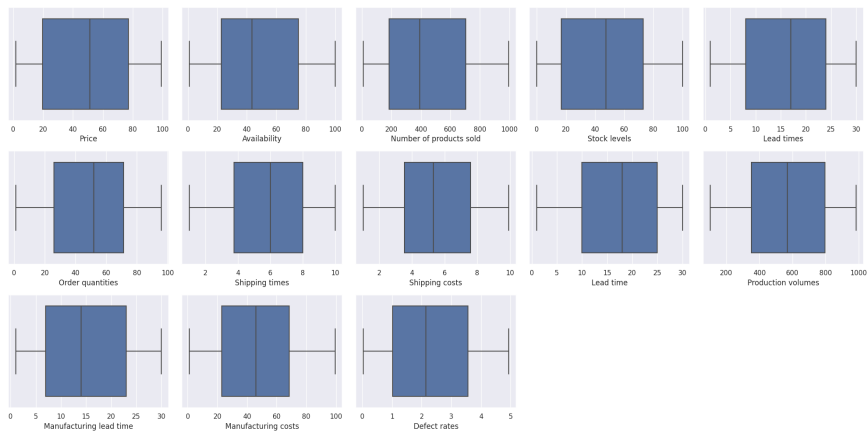


```
num_vars=['Price','Availability','Number of products sold','Stock levels','Lead times','Order quantities','Shipping times','Shipping cost']
fig,axs=plt.subplots(nrows=3,ncols=5,figsize=(20,10))
axs=axs.flatten()
```

```
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()
```

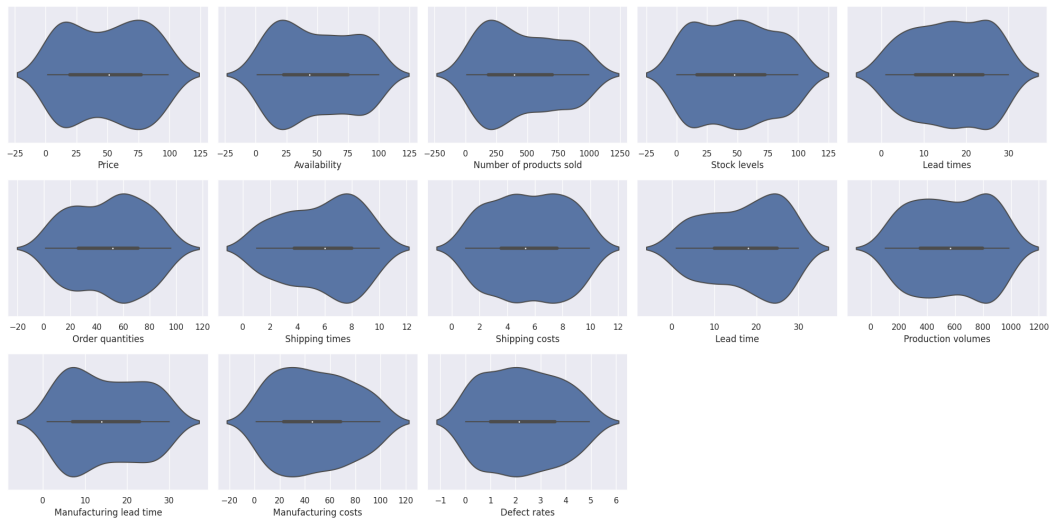


```
num_vars=['Price','Availability','Number of products sold','Stock levels','Lead times','Order quantities','Shipping times','Shipping cost']
fig,axs=plt.subplots(nrows=3,ncols=5,figsize=(20,10))
axs=axs.flatten()
```

```
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()
```

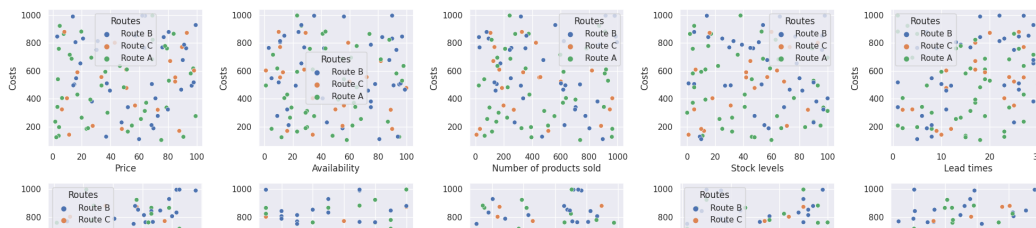


```
num_vars=['Price','Availability','Number of products sold','Stock levels','Lead times','Order quantities','Shipping times','Shipping cost']
fig,axs=plt.subplots(nrows=3,ncols=5,figsize=(20,10))
axs=axs.flatten()

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()
```

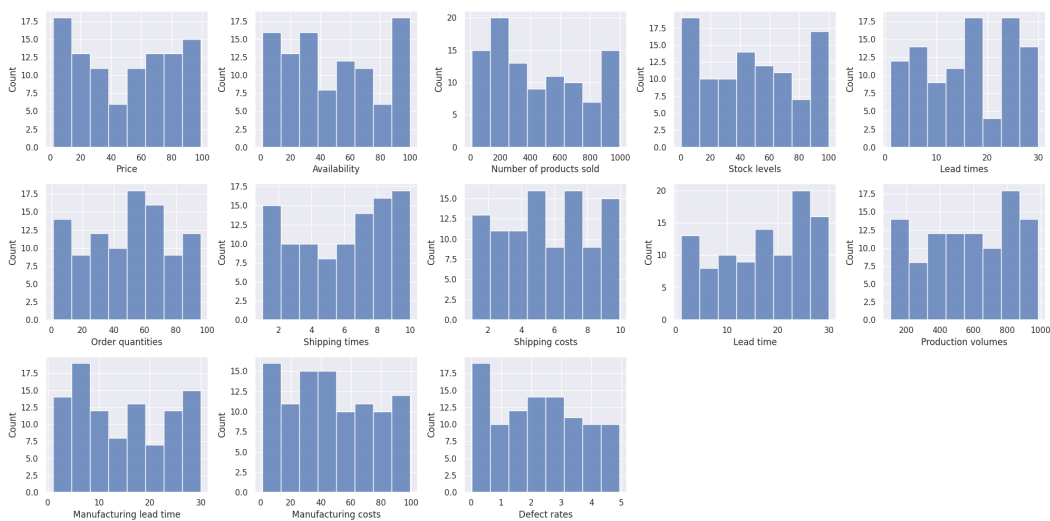


```
num_vars=['Price','Availability','Number of products sold','Stock levels','Lead times','Order quantities','Shipping times','Shipping cost']
fig,axs=plt.subplots(nrows=3,ncols=5,figsize=(20,10))
axs=axs.flatten()
```

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

```

There is no outlier so we dont have to remove it

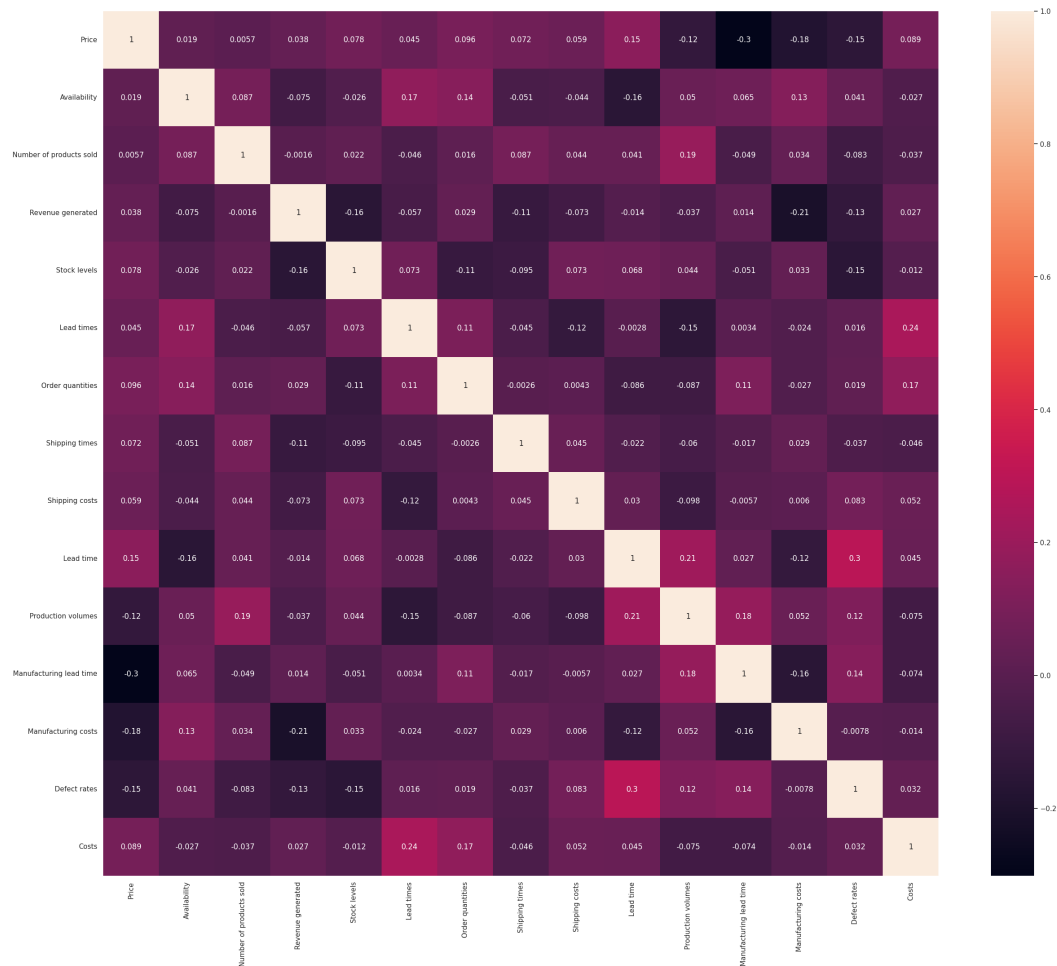
Correlation Heatmap

```

# Correlation Heatmap
plt.figure(figsize=(30,25))
sns.heatmap(df.corr(),fmt='.2g',annot=True)

```

```
<ipython-input-126-0de0d3c6ea7d>:3: FutureWarning: The default value of numeric_only in DataFrame.corr
sns.heatmap(df.corr(),fmt='.2g',annot=True)
<Axes: >
```



Train test Split

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

#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

```
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import GridSearchCV

# 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(dtrees,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_)

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

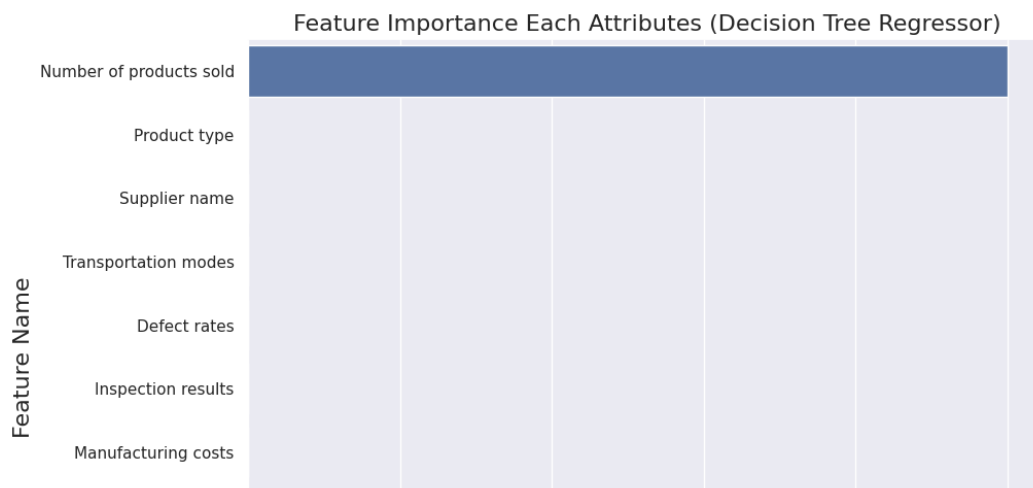
```
▼ DecisionTreeRegressor
DecisionTreeRegressor(max_depth=2, max_features='sqrt', min_samples_leaf=3,
                      random_state=0)
```

```
from sklearn import metrics
from sklearn.metrics import mean_absolute_percentage_error
import math
y_pred=dtrees.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

imp_df=pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": dtrees.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=16)
plt.xlabel("Importance",fontsize=16)
plt.ylabel("Feature Name",fontsize=16)
plt.show()
```

Random Forest Regressor

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

# create a RandomForestRegressor 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_)

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

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

```
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))

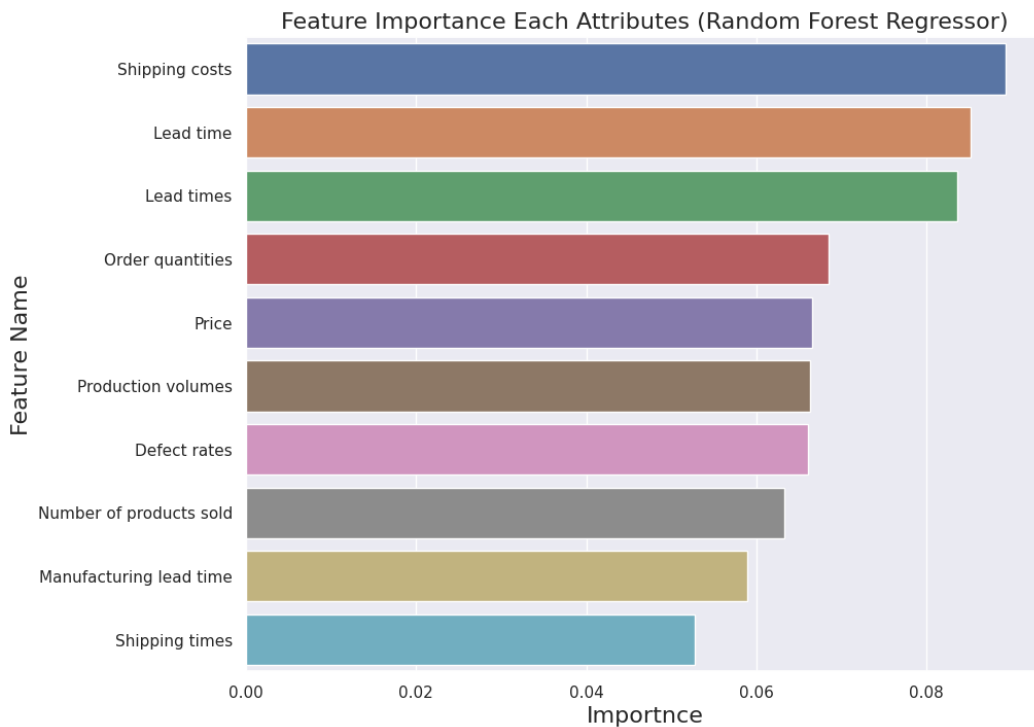
MAE is 247.33969719962744
MAPE is 0.6029768224226728
MSE is 71899.28833186119
R2 Score is -0.07294105713825938
RMSE Score is 268.14042651540103
```

```
imp_df=pd.DataFrame({
    "Feature Name": X_train.columns,
```

```

"Importance": rf.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=16)
plt.xlabel("Importance",fontsize=16)
plt.ylabel("Feature Name",fontsize=16)
plt.show()

```



AdaBoost Regressor

```

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

# create a AdaBoostRegressor 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
AdaBoostRegressor(learning_rate=1, random_state=7)

```

```

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

```

```

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=16)
plt.xlabel("Importance",fontsize=16)
plt.ylabel("Feature Name",fontsize=16)
plt.show()

```

