

IMPORTING LIBRARIES:

```
In [2]: 1 import pandas as pd
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 import numpy as np
5 from sklearn.preprocessing import StandardScaler, OneHotEncoder
6 from sklearn.pipeline import Pipeline
7 from sklearn.linear_model import LinearRegression
8 from sklearn.model_selection import train_test_split, GridSearchCV
9 from sklearn.metrics import r2_score
10 from sklearn.ensemble import RandomForestRegressor
11 import plotly.express as px
12 import warnings
13 warnings.filterwarnings("ignore")
14 from sklearn.tree import DecisionTreeRegressor
15 import xgboost as XGBRegressor
16 from sklearn.metrics import mean_squared_error
17 from sklearn.feature_selection import SelectKBest, chi2, mutual_info_classif, f_regression
18 !pip install xgboost==1.0.1
19 import xgboost as xgb
20 import matplotlib.pyplot as plt
```

Requirement already satisfied: xgboost==1.0.1 in c:\users\ruchi\anaconda3\lib\site-packages (1.0.1)

Requirement already satisfied: numpy in c:\users\ruchi\anaconda3\lib\site-packages (from xgboost==1.0.1) (1.24.3)

Requirement already satisfied: scipy in c:\users\ruchi\anaconda3\lib\site-packages (from xgboost==1.0.1) (1.10.1)

```
In [3]: 1 df=pd.read_csv("C:\\Users\\ruchi\\OneDrive\\Documents\\ML_dataset.csv")
```

```
In [4]: 1 pd.set_option("display.max_columns",100)
        2 df.head()
```

Out[4]:

| | * | Project Code | PQ # | PO / SO # | ASN/DN # | Country | Managed By | Fulfill Via | Vendor INCO Term | Shipment Mode | PQ First Sent to Client Date | PO Sent to Vendor Date | Scheduled Delivery Date | Delivered to Client Date | Delivery Recorded Date | Produ Groi |
|---|----|--------------|----------------|-----------|----------|---------------|------------|-------------|------------------|---------------|------------------------------|------------------------|-------------------------|--------------------------|------------------------|------------|
| 0 | 1 | 100-CI-T01 | Pre-PQ Process | SCMS-4 | ASN-8 | Côte d'Ivoire | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 02-Jun-06 | 02-Jun-06 | 02-Jun-06 | HRI |
| 1 | 3 | 108-VN-T01 | Pre-PQ Process | SCMS-13 | ASN-85 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 14-Nov-06 | 14-Nov-06 | 14-Nov-06 | Af |
| 2 | 4 | 100-CI-T01 | Pre-PQ Process | SCMS-20 | ASN-14 | Côte d'Ivoire | PMO - US | Direct Drop | FCA | Air | Pre-PQ Process | Date Not Captured | 27-Aug-06 | 27-Aug-06 | 27-Aug-06 | HRI |
| 3 | 15 | 108-VN-T01 | Pre-PQ Process | SCMS-78 | ASN-50 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 01-Sep-06 | 01-Sep-06 | 01-Sep-06 | Af |
| 4 | 16 | 108-VN-T01 | Pre-PQ Process | SCMS-81 | ASN-55 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 11-Aug-06 | 11-Aug-06 | 11-Aug-06 | Af |

In [5]: 1 df.describe()

Out[5]:

| | | * Unit of Measure (Per Pack) | Line Item Quantity | Line Item Value | Pack Price | Unit Price | Line Item Insurance (USD) |
|-------|--------------|------------------------------|--------------------|-----------------|--------------|--------------|---------------------------|
| count | 10324.000000 | 10324.000000 | 10324.000000 | 1.032400e+04 | 10324.000000 | 10324.000000 | 10037.000000 |
| mean | 51098.968229 | 77.990895 | 18332.534870 | 1.576506e+05 | 21.910241 | 0.611701 | 240.117626 |
| std | 31944.332496 | 76.579764 | 40035.302961 | 3.452921e+05 | 45.609223 | 3.275808 | 500.190568 |
| min | 1.000000 | 1.000000 | 1.000000 | 0.000000e+00 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 12795.750000 | 30.000000 | 408.000000 | 4.314593e+03 | 4.120000 | 0.080000 | 6.510000 |
| 50% | 57540.500000 | 60.000000 | 3000.000000 | 3.047147e+04 | 9.300000 | 0.160000 | 47.040000 |
| 75% | 83648.250000 | 90.000000 | 17039.750000 | 1.664471e+05 | 23.592500 | 0.470000 | 252.400000 |
| max | 86823.000000 | 1000.000000 | 619999.000000 | 5.951990e+06 | 1345.640000 | 238.650000 | 7708.440000 |

In [6]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 10324 entries, 0 to 10323
```

```
Data columns (total 33 columns):
```

| # | Column | Non-Null Count | Dtype |
|----|------------------------------|----------------|---------|
| 0 | * | 10324 non-null | int64 |
| 1 | Project Code | 10324 non-null | object |
| 2 | PQ # | 10324 non-null | object |
| 3 | PO / SO # | 10324 non-null | object |
| 4 | ASN/DN # | 10324 non-null | object |
| 5 | Country | 10324 non-null | object |
| 6 | Managed By | 10324 non-null | object |
| 7 | Fulfill Via | 10324 non-null | object |
| 8 | Vendor INCO Term | 10324 non-null | object |
| 9 | Shipment Mode | 9964 non-null | object |
| 10 | PQ First Sent to Client Date | 10324 non-null | object |
| 11 | PO Sent to Vendor Date | 10324 non-null | object |
| 12 | Scheduled Delivery Date | 10324 non-null | object |
| 13 | Delivered to Client Date | 10324 non-null | object |
| 14 | Delivery Recorded Date | 10324 non-null | object |
| 15 | Product Group | 10324 non-null | object |
| 16 | Sub Classification | 10324 non-null | object |
| 17 | Vendor | 10324 non-null | object |
| 18 | Item Description | 10324 non-null | object |
| 19 | Molecule/Test Type | 10324 non-null | object |
| 20 | Brand | 10324 non-null | object |
| 21 | Dosage | 8588 non-null | object |
| 22 | Dosage Form | 10324 non-null | object |
| 23 | Unit of Measure (Per Pack) | 10324 non-null | int64 |
| 24 | Line Item Quantity | 10324 non-null | int64 |
| 25 | Line Item Value | 10324 non-null | float64 |
| 26 | Pack Price | 10324 non-null | float64 |
| 27 | Unit Price | 10324 non-null | float64 |
| 28 | Manufacturing Site | 10324 non-null | object |
| 29 | First Line Designation | 10324 non-null | object |
| 30 | Weight (Kilograms) | 10324 non-null | object |
| 31 | Freight Cost (USD) | 10324 non-null | object |
| 32 | Line Item Insurance (USD) | 10037 non-null | float64 |

```
dtypes: float64(4), int64(3), object(26)
```

```
memory usage: 2.6+ MB
```

DATA PRE-PROCESSING:

```
In [7]: 1 date_list=['Scheduled Delivery Date','Delivered to Client Date','Delivery Recorded Date']
        2 same_date=[]
        3 for i in range(len(date_list)-1):
        4     for j in df.index:
        5         if df[date_list[i]][j]!=df[date_list[i+1]][j]:
        6             same_date.append(j)
```

```
In [8]: 1 len(same_date)
```

Out[8]: 5911

```
In [9]: 1 date_list=['Scheduled Delivery Date','Delivered to Client Date','Delivery Recorded Date']
        2 df['Date_deliver_client']=np.nan
        3 df['Date_Delivery']=np.nan
        4 df['Date_deliver_record']=np.nan
        5
        6 target=['Date_deliver_client','Date_Delivery','Date_deliver_record']
        7 def date_input(target,present):
        8     for j in df.index:
        9         df[target][j]= df[present][j].split('-')[0]
        10
        11 for i in range(0,3):
        12     date_input(target[i],date_list[i])
```


```
In [10]: 1 df['month_deliver_client']=np.nan
        2 df['month_Delivery']=np.nan
        3 df['month_deliver_record']=np.nan
        4 target=['month_deliver_client','month_Delivery','month_deliver_record']
        5 def month_input(target,present):
        6     for j in df.index:
        7         df[target][j]= df[present][j].split('-')[1]
        8 for i in range(0,3):
        9     month_input(target[i],date_list[i])
```

```
In [11]: 1 df['year_deliver_client']=np.nan
2 df['year_Delivery']=np.nan
3 df['year_deliver_record']=np.nan
4 target=['year_deliver_client','year_Delivery','year_deliver_record']
5 def year_input(target,present):
6     for j in df.index:
7         df[target][j]= df[present][j].split('-')[2]
8 for i in range(0,3):
9     year_input(target[i],date_list[i])
```

```
In [12]: 1 df.head()
```

Out[12]:

| | * | Project Code | PQ # | PO / SO # | ASN/DN # | Country | Managed By | Fulfill Via | Vendor INCO Term | Shipment Mode | PQ First Sent to Client Date | PO Sent to Vendor Date | Scheduled Delivery Date | Delivered to Client Date | Delivery Recorded Date | Produ Groi |
|---|----|--------------|----------------|-----------|----------|---------------|------------|-------------|------------------|---------------|------------------------------|------------------------|-------------------------|--------------------------|------------------------|------------|
| 0 | 1 | 100-CI-T01 | Pre-PQ Process | SCMS-4 | ASN-8 | Côte d'Ivoire | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 02-Jun-06 | 02-Jun-06 | 02-Jun-06 | HRI |
| 1 | 3 | 108-VN-T01 | Pre-PQ Process | SCMS-13 | ASN-85 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 14-Nov-06 | 14-Nov-06 | 14-Nov-06 | Af |
| 2 | 4 | 100-CI-T01 | Pre-PQ Process | SCMS-20 | ASN-14 | Côte d'Ivoire | PMO - US | Direct Drop | FCA | Air | Pre-PQ Process | Date Not Captured | 27-Aug-06 | 27-Aug-06 | 27-Aug-06 | HRI |
| 3 | 15 | 108-VN-T01 | Pre-PQ Process | SCMS-78 | ASN-50 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 01-Sep-06 | 01-Sep-06 | 01-Sep-06 | Af |
| 4 | 16 | 108-VN-T01 | Pre-PQ Process | SCMS-81 | ASN-55 | Vietnam | PMO - US | Direct Drop | EXW | Air | Pre-PQ Process | Date Not Captured | 11-Aug-06 | 11-Aug-06 | 11-Aug-06 | Af |



```
In [13]: 1 df.drop(date_list,axis=1,inplace=True)
```


In [14]:

```
1 df.dtypes
```

```

Out[14]: *
Project Code      int64
PQ #             object
PO / SO #        object
ASN/DN #         object
Country          object
Managed By      object
Fulfill Via      object
Vendor INCO Term object
Shipment Mode    object
PQ First Sent to Client Date object
PO Sent to Vendor Date object
Product Group    object
Sub Classification object
Vendor           object
Item Description  object
Molecule/Test Type object
Brand            object
Dosage           object
Dosage Form      object
Unit of Measure (Per Pack) int64
Line Item Quantity int64
Line Item Value   float64
Pack Price        float64
Unit Price        float64
Manufacturing Site object
First Line Designation object
Weight (Kilograms) object
Freight Cost (USD) object
Line Item Insurance (USD) float64
Date_deliver_client object
Date_Delivery     object
Date_deliver_record object
month_deliver_client object
month_Delivery    object
month_deliver_record object
year_deliver_client object
year_Delivery     object
year_deliver_record object
dtype: object

```

```
In [15]: 1 wt_list=df['Weight (Kilograms)'].tolist()
2 wt_new={}
3 for i in range(len(wt_list)):
4     try:
5         pd.to_numeric(wt_list[i])
6
7     except:
8         wt_new[i]=wt_list[i]
9 len(wt_new)
10 df['wt_cat']=np.nan
```

```
In [16]: 1 for i in wt_new.keys():
2     df['wt_cat'][i]=wt_new[i]
3     df['Weight (Kilograms)'][i]=0
```

```
In [17]: 1 df['Weight (Kilograms)']=df['Weight (Kilograms)'].astype(int)
```

```
In [18]: 1 ft_list=df['Freight Cost (USD)'].tolist()
2 ft_new={}
3 for i in range(len(ft_list)):
4     try:
5         pd.to_numeric(ft_list[i])
6
7     except:
8         ft_new[i]=ft_list[i]
9 len(ft_new)
10 df['ft_cat']=np.nan
11 for i in ft_new.keys():
12     df['ft_cat'][i]=ft_new[i]
13     df['Freight Cost (USD)'][i]=0
14 df['Freight Cost (USD)']=df['Freight Cost (USD)'].astype(float)
```

In [19]:

```
1 df.dtypes
```

```
Out[19]: *
Project Code      int64
PQ #              object
PO / SO #         object
ASN/DN #          object
Country           object
Managed By       object
Fulfill Via       object
Vendor INCO Term  object
Shipment Mode     object
PQ First Sent to Client Date object
PO Sent to Vendor Date object
Product Group     object
Sub Classification object
Vendor            object
Item Description   object
Molecule/Test Type object
Brand             object
Dosage            object
Dosage Form       object
Unit of Measure (Per Pack) int64
Line Item Quantity int64
Line Item Value    float64
Pack Price         float64
Unit Price         float64
Manufacturing Site object
First Line Designation object
Weight (Kilograms) int32
Freight Cost (USD) float64
Line Item Insurance (USD) float64
Date_deliver_client object
Date_Delivery      object
Date_deliver_record object
month_deliver_client object
month_Delivery     object
month_deliver_record object
year_deliver_client object
year_Delivery      object
year_deliver_record object
wt_cat            object
```

```
ft_cat          object
dtype: object
```

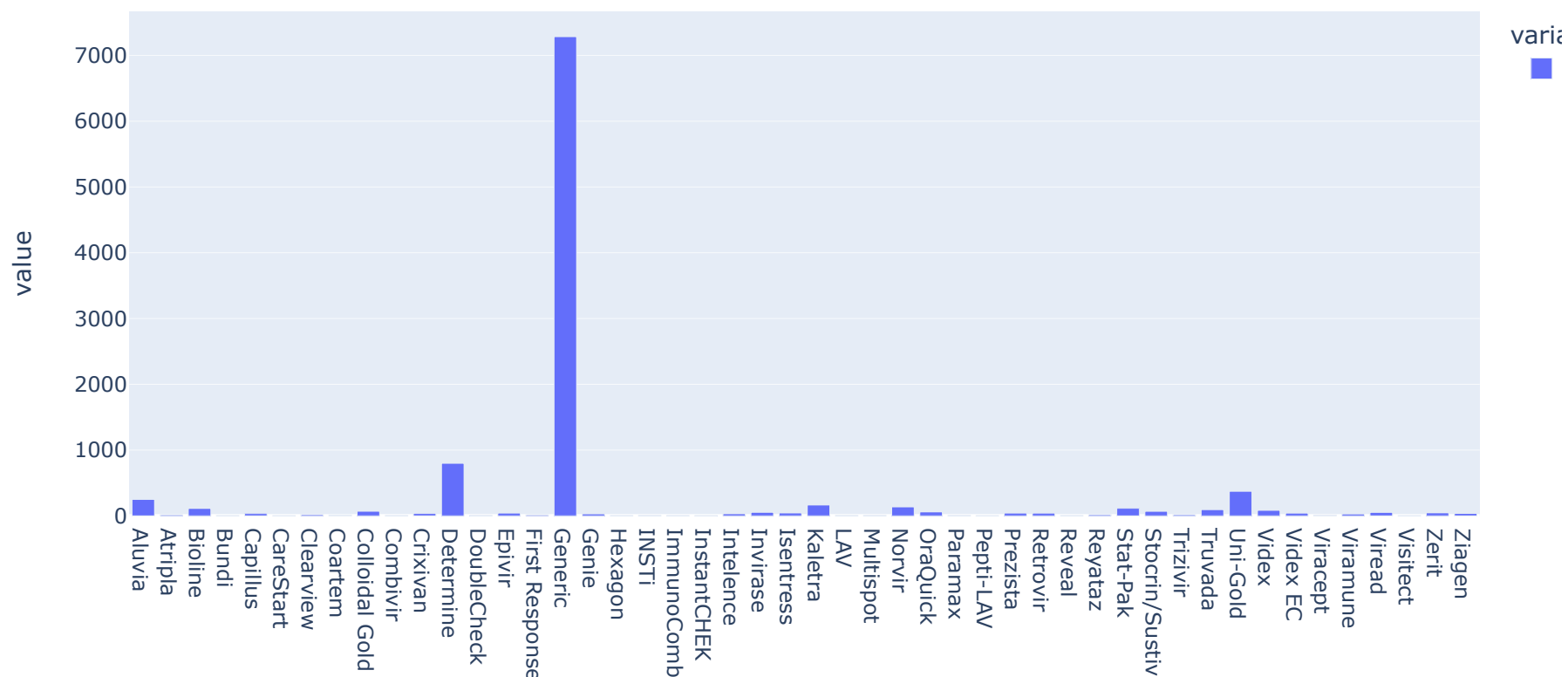
```
In [20]: 1 df['Dosage Form'].fillna(-1,inplace=True)
          2 df['Line Item Insurance (USD)'].fillna(0,inplace=True)
          3 df.to_excel("C:final_sheet.xlsx")
```

DATA VISUALIZATION:

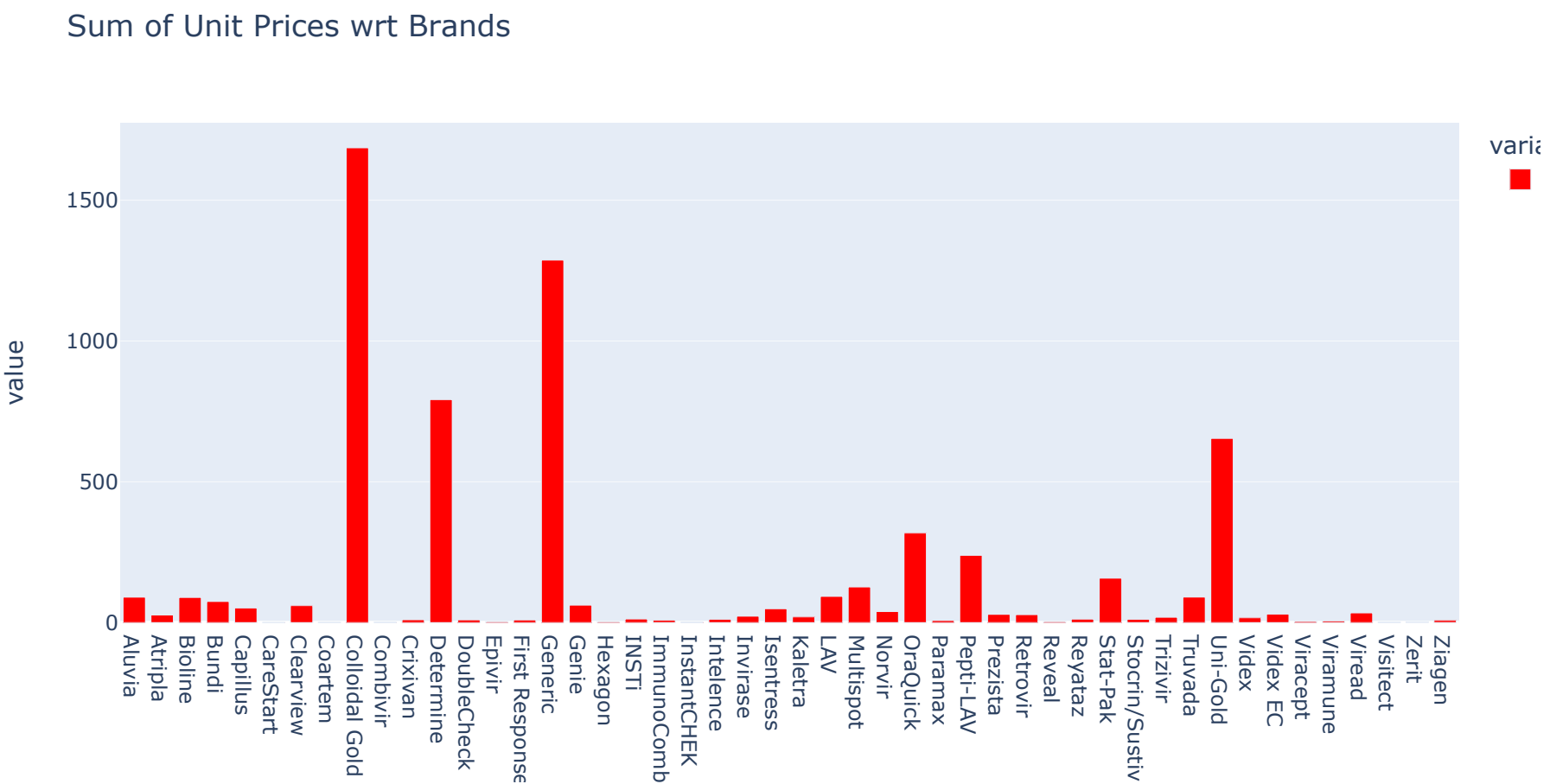
```
In [21]: 1 print("The total sum of unit price is ", df[df['Brand']=='Generic']['Unit Price'].sum())
2 px.bar(df.groupby('Brand')['Unit Price'].count(), title='Count of Unit Prices wrt Brands')
```

The total sum of unit price is 1287.3200000000002

Count of Unit Prices wrt Brands

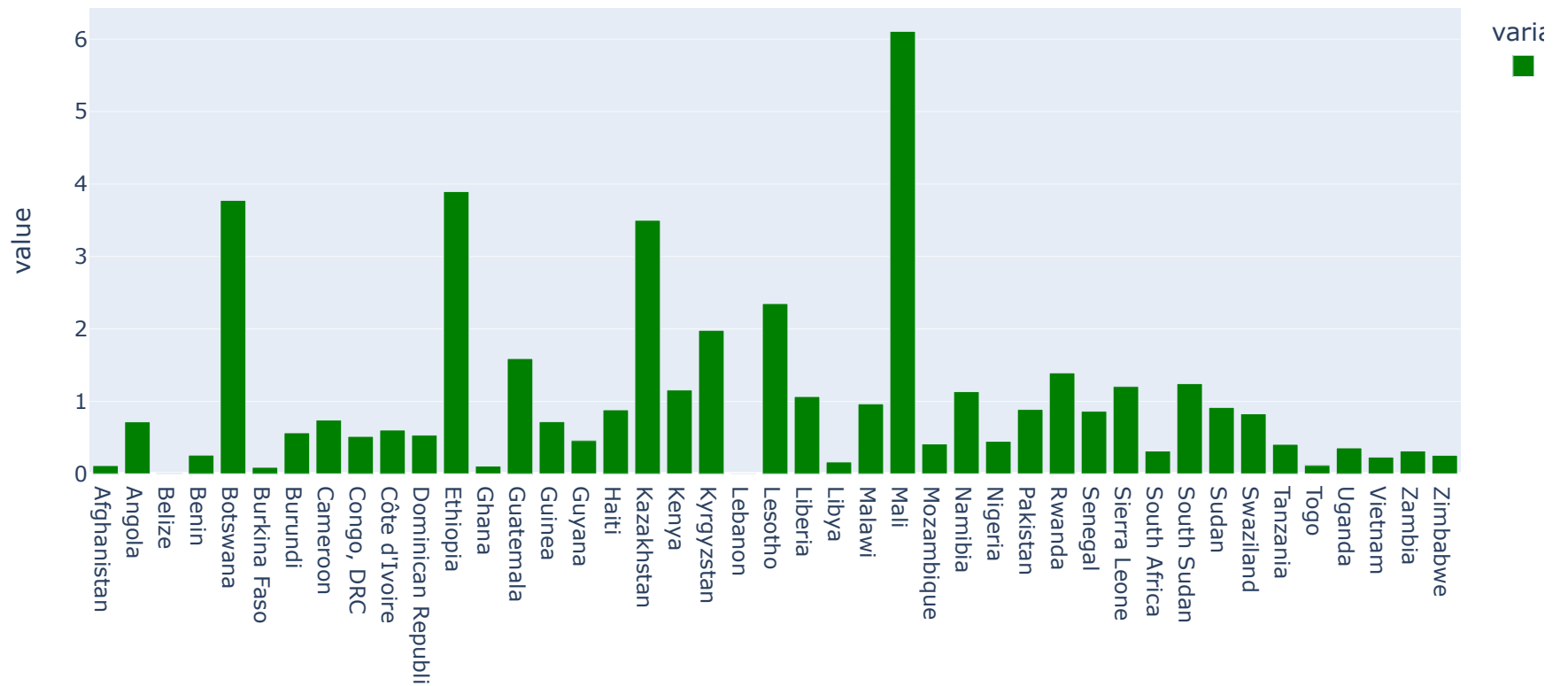


```
In [22]: 1 px.bar(df.groupby('Brand')['Unit Price'].sum(), title='Sum of Unit Prices wrt Brands',color_discrete_sequence = ['r
```



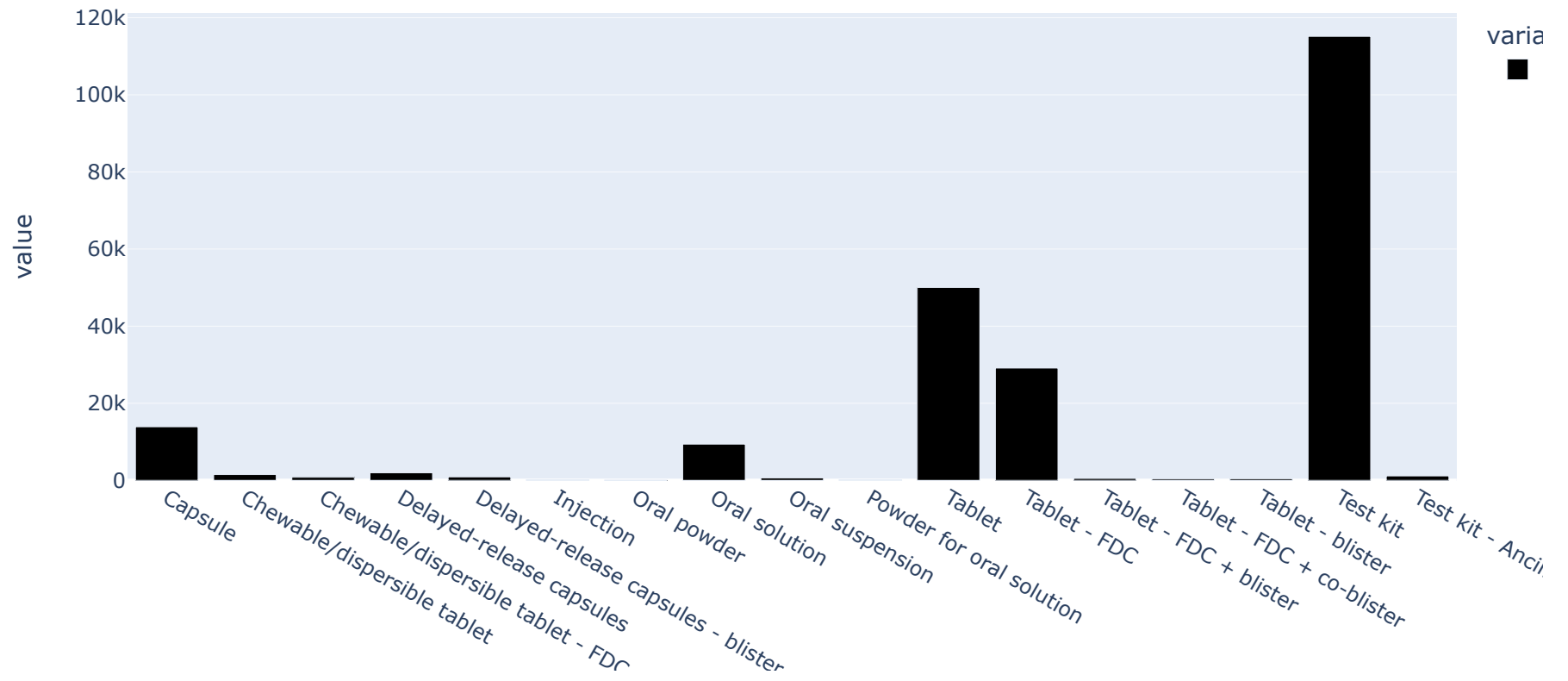

```
In [23]: 1 px.bar(df.groupby('Country')['Unit Price'].sum()/df.groupby('Country')['Unit Price'].count(), title='Ratio of Sum
2
```

Ratio of Sum to Count of Unit Prices wrt Country



```
In [24]: 1 px.bar(df.groupby('Dosage Form')['Pack Price'].sum(), title='Sum of Pack Prices wrt Vendor',color_discrete_sequenc
```

Sum of Pack Prices wrt Vendor



INITIALIZATION:


```
In [25]: 1 X=df.drop(['Pack Price','Unit Price'],axis=1)
```

```
In [26]: 1 y=df['Pack Price']  
2 Y=df['Unit Price']
```

```
In [27]: 1 cat_col=[]  
2 num_col=[]  
3 for i in X.columns:  
4     if X[i].dtype=='O':  
5         cat_col.append(i)  
6     else:  
7         num_col.append(i)  
8 categorical_transformer = OneHotEncoder(handle_unknown='ignore')
```

```
In [28]: 1 import pickle  
2 pickle.dump(categorical_transformer, open('cattrans.pkl','wb'))
```

```
In [29]: 1 from sklearn.preprocessing import MinMaxScaler  
2 from sklearn.pipeline import Pipeline  
3  
4 numeric_transformer = Pipeline(steps=[('scaler',MinMaxScaler()) ])  
5 pickle.dump(numeric_transformer, open('numerictrans.pkl','wb'))  
6 preprocessing=ColumnTransformer(remainder='passthrough', transformers=[('cat', categorical_transformer, cat_col),  
7 pickle.dump(preprocessing, open('prepro.pkl','wb'))
```



#FOR UNITPRICE:

```
In [30]: 1 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2,random_state=0)
```

In [31]:

```
1 import pandas as pd
2
3 date_str = '12-05-10'
4 date_obj = pd.to_datetime(date_str, format='%y-%m-%d')
5 timestamp = date_obj.timestamp() # Converts to a float representing the timestamp
6 # Example usage
7 print(timestamp)
8
9 final_list=[]
10 model=[LinearRegression(),DecisionTreeRegressor(),RandomForestRegressor(),xgb.XGBRegressor()]
11 result={}
12 for j in [5, 10, 15, 20, '25', '30', '35', '39']:
13     new=[]
14     for i in model:
15
16         runpipe = Pipeline(steps=[('preprocessor', preprocessing), ('skb',SelectKBest(score_func=f_regression,
17                                     ('model', i)))]
18         runpipe.fit(X_train,Y_train)
19         pred=runpipe.predict(X_test)
20         rmse = np.sqrt(mean_squared_error(Y_test, pred))
21         new.append(runpipe.score(X_test, Y_test))
22     final_list.append(new)
```

1336608000.0

```
In [32]: 1 pd.DataFrame(final_list,index=['5', '10', '15', '20', '25', '30','35','39'], columns=['Linear Regression','Decision
```

Out[32]:

| | Linear Regression | Decision Tree Regressor | Random Forest Regressor | XGBRegressor() |
|-----------|-------------------|-------------------------|-------------------------|----------------|
| 5 | -0.002967 | -0.002967 | -0.002678 | -0.002967 |
| 10 | 0.731158 | 0.731158 | 0.731148 | 0.731158 |
| 15 | 0.816197 | 0.816197 | 0.815141 | 0.816198 |
| 20 | 0.820665 | 0.824756 | 0.825527 | 0.824757 |
| 25 | 0.820532 | 0.824191 | 0.823992 | 0.824193 |
| 30 | 0.821485 | 0.825144 | 0.824633 | 0.825166 |
| 35 | 0.891285 | 0.895023 | 0.895634 | 0.895026 |
| 39 | 0.890749 | 0.852023 | 0.852094 | 0.851924 |

#FOR PACK_PRICE:

```

In [33]: 1 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2,
2                                                    random_state=0)
3 final_list=[]
4 model=[LinearRegression(),DecisionTreeRegressor(),RandomForestRegressor(),xgb.XGBRegressor()]
5 result={}
6 for j in [5, 10, '15', '20', '25', '30','35','39']:
7     new=[]
8     for i in model:
9
10         runpipe = Pipeline(steps=[('preprocessor', preprocessing), ('skb',SelectKBest(score_func=f_regression,
11                                     ('model', i)))]
12         runpipe.fit(X_train,y_train)
13         pred=runpipe.predict(X_test)
14         rmse = np.sqrt(mean_squared_error(y_test, pred))
15         new.append(runpipe.score(X_test, y_test))
16     final_list.append(new)
17 pd.DataFrame(final_list,index=['5', '10', '15', '20', '25', '30','35','39'], columns=['Linear Regression','Decision

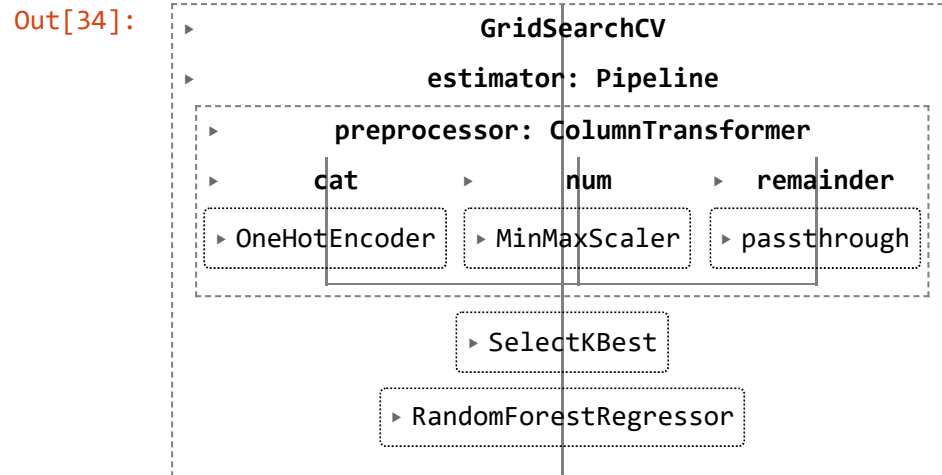
```

Out[33]:

| | Linear Regression | Decision Tree Regressor | Random Forest Regressor | XGBRegressor() |
|-----------|-------------------|-------------------------|-------------------------|----------------|
| 5 | -0.002967 | -0.002967 | -0.003079 | -0.002967 |
| 10 | 0.731158 | 0.731158 | 0.731174 | 0.731158 |
| 15 | 0.816197 | 0.816197 | 0.816519 | 0.816198 |
| 20 | 0.820665 | 0.824756 | 0.824899 | 0.824757 |
| 25 | 0.820532 | 0.824191 | 0.825693 | 0.824193 |
| 30 | 0.821485 | 0.825144 | 0.825129 | 0.825166 |
| 35 | 0.891285 | 0.895023 | 0.894615 | 0.895026 |
| 39 | 0.890749 | 0.852023 | 0.853343 | 0.851924 |

HYPERTUNING OF RANDOM FOREST

```
In [34]: 1 pipeline = Pipeline(steps=[('preprocessor', preprocessing), ('skb', SelectKBest(score_func=f_regression, k = int(j)
2                                     ('model', RandomForestRegressor()))])
3 parameters={'model__max_depth':[50, None],
4             'model__min_samples_split':[2, 5],
5             'model__min_samples_leaf':[1, 2]}
6 pickle.dump(preprocessing, open('prepro.pkl', 'wb'))
7 CV = GridSearchCV(pipeline, parameters, scoring = 'neg_mean_absolute_error', n_jobs= 1)
8 CV.fit(X_train, y_train)
```



```
In [35]: 1 CV.best_params_
```

```
Out[35]: {'model__max_depth': 50,
          'model__min_samples_leaf': 1,
          'model__min_samples_split': 2}
```

```
In [36]: 1 y_pred = CV.predict(X_test)
2 r2_score(y_test, y_pred)
3
```

```
Out[36]: 0.8545084237273932
```

In []: 1