

## sales.py

```
1 # Importing necessary libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.express as px
7 from sklearn.ensemble import RandomForestRegressor
8 from sklearn.preprocessing import LabelEncoder
9 from sklearn.model_selection import train_test_split
10
11 # Load data
12 train = pd.read_csv("../input/tabular-playground-series-jan-2022/train.csv",
13                     index_col="row_id")
14 test = pd.read_csv("../input/tabular-playground-series-jan-2022/test.csv",
15                    index_col="row_id")
16 sample = pd.read_csv("../input/tabular-playground-series-jan-2022/sample_submission.csv")
17
18 # Convert 'date' column to datetime format
19 train['date'] = pd.to_datetime(train['date'])
20 test['date'] = pd.to_datetime(test['date'])
21
22 # EDA and feature engineering
23 # Add extra features
24 train['year'] = train['date'].dt.year
25 train['month'] = train['date'].dt.month_name()
26 train['day'] = train['date'].dt.day_name()
27
28 test['year'] = test['date'].dt.year
29 test['month'] = test['date'].dt.month_name()
30 test['day'] = test['date'].dt.day_name()
31
32 # Add weekend feature
33 train['is_weekend'] = train['day'].apply(lambda x: 1 if x in ['Saturday', 'Sunday'] else 0)
34 test['is_weekend'] = test['day'].apply(lambda x: 1 if x in ['Saturday', 'Sunday'] else 0)
35
36 # Add time step feature
37 train['Time_step'] = np.arange(len(train))
38 test['Time_step'] = np.arange(len(test))
39
40 # Label encoding for categorical features
41 le = LabelEncoder()
42 cols = ['country', 'product', 'store', 'month', 'day', 'year']
43 for col in cols:
44     train[col] = le.fit_transform(train[col])
45     test[col] = le.transform(test[col])
46
47 # Train-validation split
48 X_train = train[train.date <= '2018-05-31'].drop(['date', 'num_sold'], axis=1)
49 X_val = train[(train.date >= '2018-06-01') & (train.date <= '2018-12-31')].drop(['date',
50 'num_sold'], axis=1)
51 y_train = train[train.date <= '2018-05-31']['num_sold']
52 y_val = train[(train.date >= '2018-06-01') & (train.date <= '2018-12-31')]['num_sold']
```

```

50
51 # Model - Random Forest Regressor
52 baseline_regressor = RandomForestRegressor(n_estimators=500, n_jobs=-1)
53 baseline_regressor.fit(X_train, y_train)
54
55 # Validation predictions
56 val_pred = baseline_regressor.predict(X_val)
57
58 # SMAPE evaluation function
59 def SMAPE(y_true, y_pred):
60     denominator = (y_true + np.abs(y_pred)) / 200.0
61     diff = np.abs(y_true - y_pred) / denominator
62     diff[denominator == 0] = 0.0
63     return np.mean(diff)
64
65 # Calculate SMAPE
66 print('The SMAPE value of the Random Forests model is:', SMAPE(y_val, val_pred))
67
68 # Feature importance
69 importance = baseline_regressor.feature_importances_
70 feature_names = list(X_train.columns)
71 std = np.std([tree.feature_importances_ for tree in baseline_regressor.estimators_], axis=0)
72 baseline_importances = pd.Series(importance, index=feature_names)
73
74 # Plot feature importance
75 fig, ax = plt.subplots()
76 baseline_importances.plot.bar(xerr=std, ax=ax)
77 ax.set_title("Feature Importances")
78 ax.set_ylabel("Mean decrease in impurity")
79 fig.tight_layout()
80
81 # Predictions for test data
82 test_pred = baseline_regressor.predict(test.drop(['date'], axis=1))
83
84 # Prepare submission
85 sample['num_sold'] = test_pred
86 sample.to_csv("baseliner.csv", index=False)
87
88 # Visualization of predicted sales trends
89 test['num_sold'] = test_pred
90 data = pd.concat([train, test])
91
92 data_plot = data.groupby(['date']).mean().reset_index()
93 plt.figure(figsize=(15, 7))
94 sns.lineplot(x=data_plot.date, y=data_plot.num_sold)
95 plt.title('Number Sold Over Time')
96 plt.show()
97 import pandas as pd
98 import numpy as np
99 from statsmodels.tsa.arima.model import ARIMA
100 from sklearn.model_selection import train_test_split
101 from sklearn.ensemble import RandomForestRegressor
102 from sklearn.metrics import mean_absolute_error, mean_squared_error
103

```

```

104 # Load Dataset
105 # Replace 'sales_data.csv' with your dataset file.
106 data = pd.read_csv('sales_data.csv', parse_dates=['Date'], index_col='Date')
107
108 # Data Preprocessing
109 data = data.asfreq('D') # Ensure daily frequency
110 data.fillna(method='ffill', inplace=True) # Fill missing values
111
112 # Visualization (optional)
113 import matplotlib.pyplot as plt
114 data['Sales'].plot(title="Sales Over Time")
115 plt.show()
116
117 # --- ARIMA Model ---
118 # Train-Test Split for Time Series
119 train_size = int(len(data) * 0.8)
120 train, test = data[:train_size], data[train_size:]
121
122 # Fit ARIMA Model
123 arima_model = ARIMA(train['Sales'], order=(5, 1, 0)) # Change parameters as needed
124 arima_fit = arima_model.fit()
125
126 # Forecast
127 arima_forecast = arima_fit.forecast(steps=len(test))
128 arima_results = pd.DataFrame({'Actual': test['Sales'], 'Forecast': arima_forecast})
129 print(arima_results)
130
131 # Evaluation Metrics for ARIMA
132 arima_mae = mean_absolute_error(test['Sales'], arima_forecast)
133 arima_rmse = np.sqrt(mean_squared_error(test['Sales'], arima_forecast))
134 print(f"ARIMA - MAE: {arima_mae}, RMSE: {arima_rmse}")
135
136 # --- Random Forest Regressor ---
137 # Feature Engineering
138 data['Day'] = data.index.day
139 data['Month'] = data.index.month
140 data['Year'] = data.index.year
141 data['Lag1'] = data['Sales'].shift(1)
142 data.dropna(inplace=True)
143
144 X = data[['Day', 'Month', 'Year', 'Lag1']]
145 y = data['Sales']
146
147 # Train-Test Split
148 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
149
150 # Fit Random Forest Model
151 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
152 rf_model.fit(X_train, y_train)
153
154 # Forecast
155 rf_forecast = rf_model.predict(X_test)
156
157 # Evaluation Metrics for Random Forest

```

```
158 rf_mae = mean_absolute_error(y_test, rf_forecast)
159 rf_rmse = np.sqrt(mean_squared_error(y_test, rf_forecast))
160 print(f"Random Forest - MAE: {rf_mae}, RMSE: {rf_rmse}")
161
162 # Visualization of Forecasts
163 plt.figure(figsize=(12, 6))
164 plt.plot(test.index, test['Sales'], label="Actual Sales")
165 plt.plot(test.index, arima_forecast, label="ARIMA Forecast")
166 plt.plot(y_test.index, rf_forecast, label="RF Forecast", linestyle='dashed')
167 plt.legend()
168 plt.title("Sales Forecasting")
169 plt.show()
```