# Research Associate KTP with the University of Essex and Croud Inc
## Technical Task

Ashkan Lotfipoor

School of Energy, Geoscience, Infrastructure and Society (EGIS),
Heriot-Watt University

July 26, 2023

# Introduction

This report is the result of my analysis of the technical assignment for the KTP Associate with the University of Essex and Croud Inc. This report is divided into two sections. The first section presents the results of the marketing data modelling on the provided dataset. Then, section two provide results for the question in Bayesian methods section of the task.

# 1    Classical Marketing Data Modelling

## 1.1    Modelling

In this section several regression models are built to determine the effectiveness of marketing strategies of the company and also to predict the revenue. Firstly the data is divided to two dataset, the train and test set, to effectively validate the results of the model on unseen data. The 2016/05/30 date is used as the split point in the dataset.

In this section, the temporal aspect of the data is not considered in the modelling process. So, the week number and Date column are dropped from the dataframe. The remaining variables or predictors in the data are media1_S (spend on media 1), media2_S (spend on media 2), media3_S (spend on media 3), competitor_sales (competitors sales), and newsletter (number of newsletter subscription). There are various statistical methods for feature selection in a regression task, such as Chi-square Test, Variance Threshold, Dispersion Ratio, Random Forest Importance and etc. I used Correlation Coefficient to select predictors. Figure 1 presents the correlation matrix. Here, the variables with Coefficient less than 0.4 were dropped from the dataframe. In the final dataset only media1_S, competitor_sales, and newsletter are used to predict revenue.
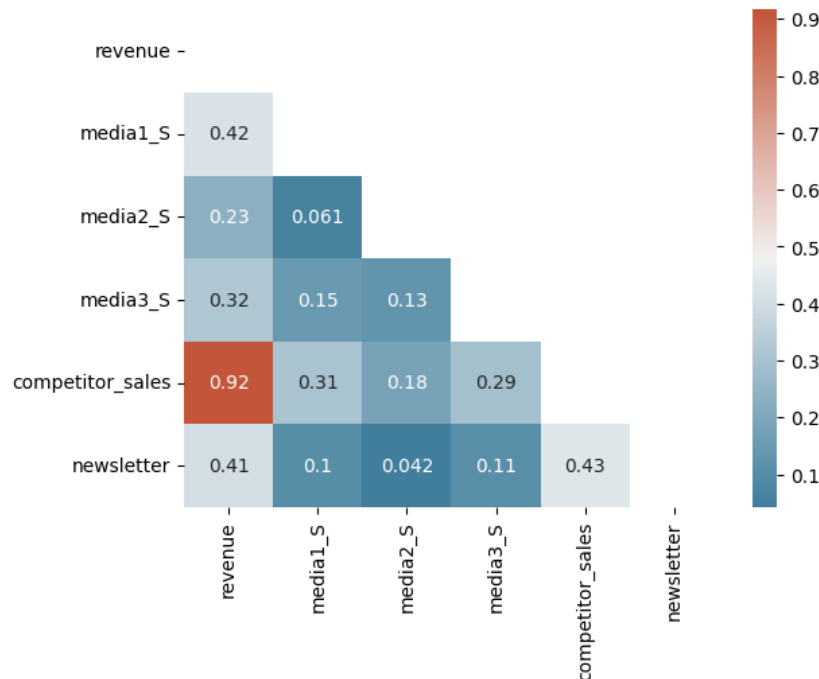


Figure 1: Correlation Matrix

For selection of the model, several commonly used regression model are used to model the marketing dataset. Using SKlearn library in Python, five algorithm are used for modelling the marketing data: linear regression, stochastic gradient descent regression, Random Forest, gradient boosting regression, Multilayer Perceptron (MLP). The results are presented in Table 1 and Figure 2.

Linear Regression and Stochastic Gradient Descent Regression models perform rather well in this task. The reported $R^2$ value for both models are 95%. Other methods does not perform well, but it is worth noting that each of these algorithms need hyperparameter tuning to perform well. Here, grid search was performed, but to reach the maximum potential of each model more analysis needs to be done. Analysing the effect of variables on

Table 1: Forecasting Error for Methods Used

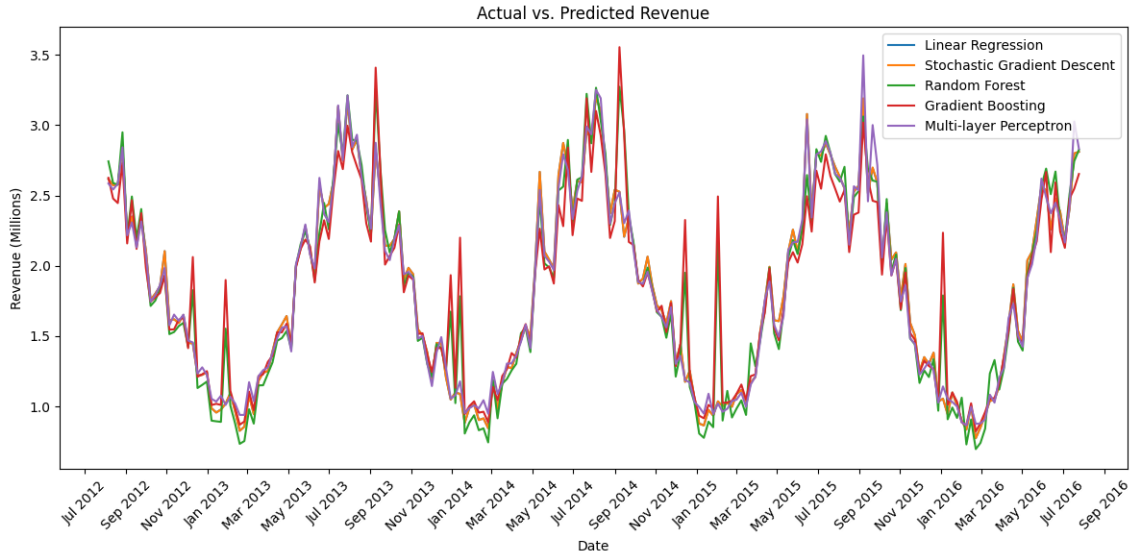| Model | RMSE | MAE | R2 | MAPE |
|---|---|---|---|---|
| Linear Regression | 41441.49 | 34845.88 | 0.952 | 1.36 |
| Stochastic Gradient Descent | 41650.05 | 34783.47 | 0.952 | 1.36 |
| Multi-layer Perceptron | 123783.20 | 90774.79 | 0.576 | 3.55 |
| Random Forest | 139584.86 | 116637.73 | 0.461 | 4.79 |
| Gradient Boosting | 145422.27 | 136118.76 | 0.415 | 5.47 |



Figure 2: Forecasted Revenue

the regression results is an essential step in understanding the relationship between the independent variables (features) and the dependent variable (target) in a regression model. Somme common methods are coefficient estimates, hypothesis tests, and residual analysis which will be performed in the next section. These are the coefficient of the regression model: 112513.37, 611400.75, and 12931.45.

The biggest limitation of the model used in the context of media marketing is the lack of model interpretability. Advanced regression models, like neural networks, can be complex and lack interpretability, making it challenging for marketers to understand the model's decision-making process. That's why the methods such as Bayesian model estimation can be usefull in this field.

## 1.2 Prediction

This section overlaps with the previous section, the data was split into two datasets and appropriate variables were used for developing the prediction model. Figure 3 presents the forecasted revenue for the test period.
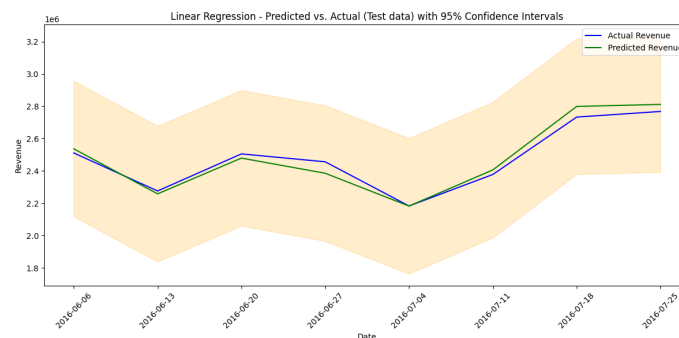


Figure 3: Forecast Results for Test Set

The performance of the model was evaluated using various metrics, RMSE, MAE, $R\hat{2}$ and MAPE. The results for the test set are presented in Table 1.
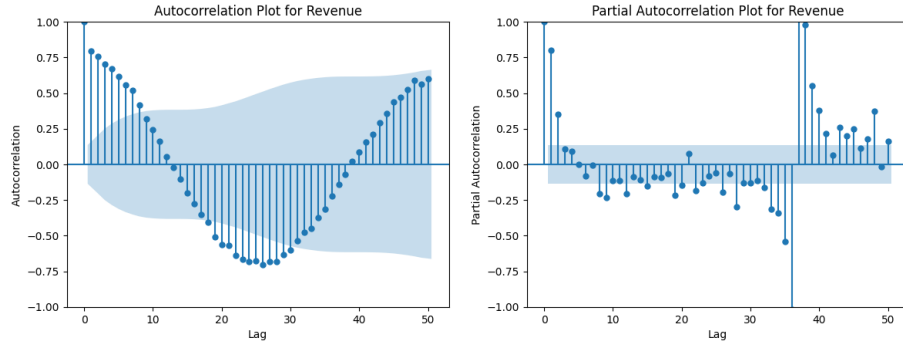
## 1.3    Temporal effects



Figure 4: ACF and PACF Plots

As it is suggested by a colleague in the task, there are a clear affects of time seen in the marketing dataset. This can be seen by plotting the revenue against time. However, to be more precise, we can plot the Autocorrelation Function (ACF) of the revenue to determine whether trends and seasonal patterns are present. Figure 4 presents the ACF plot. As the autocorrelation is significant for number of lags, we can conclude that there is a temporal effect in the dataset. From the plot is can also be concluded that there is a clear trend and seasonality in the dataset. We can investigate this further by decomposition of the dataset. Figure 5 shows the result of the time series decomposition of the revenue. This enables us to identify patterns and trends that might impact the revenue data.



Figure 5: Time Series Decomposition of the Revenue

Next, we are going to model the dataset with the temporal components of the data. For this I used ARIMA to forecast the revenue. For this, we need to first check to see if the data is stationary, as this is one the requirements of the ARIMA algorithm. If the data is not stationary, we need to perform transformation such as differencing on the revenue. For selecting the autoregressive term or the number of lag observations in

Figure 6: Line Plot and Denesity Plot of Residual Error

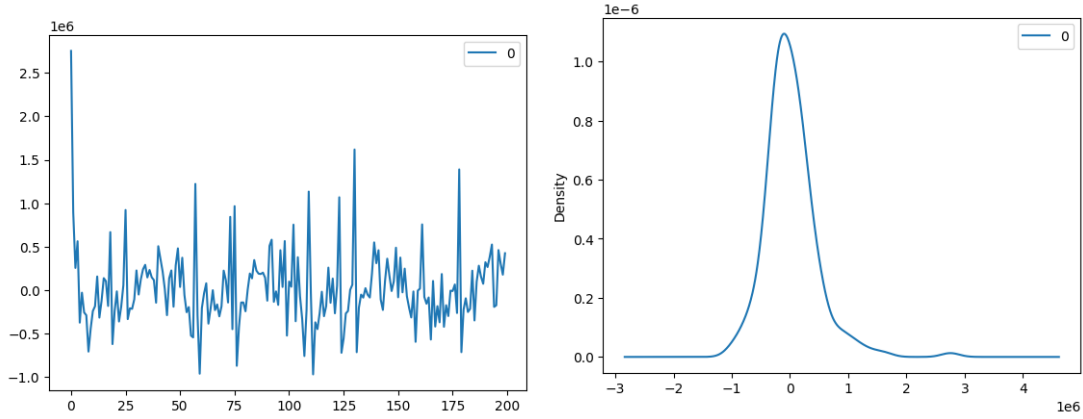the ARIMA we need to plot the Partial Autocorrelation function (PACF) which is presented in Figure 4. By analysing the ACF and PACF the number of lag was selected as 12. This value was also selected during the Grid Search process for hyperparameter optimisation. Details of this are presented in the code. The best ARIMA order (p, d, q) is selected as following: 12, 1, 2. For this the RMSE value was reported as 217787.24.

Aside from the tipycal metrics that are used for evaluating the regression error, we can also plot the residuals as a line graph and also plot the density of residuals to check if the model has captured all of patterns in the data. This is presented in Figure 6.

Using ARIMA, we considered the temporal aspect of the data, but we discard the other features in the data. We can use other algorithms to predict revenue using all of the available features along side the temporal features. This can be done by using feature engineering, we can add new features to the data, like number of week, month, and year to do this. However, using windowing method we can transform the data and keep the temporal aspect of the data too. I wrote a function in python to do this. After this we can use any algorithm to perform prediction.

In the code submitted with this report a deep learning model is used to predict the revenue. The RMSE error is 213412.58, but please bear in mind that deep learning needs a large quantity of data for the algorithm to perform. Here we only used 208 data points. Also, by spending more time on the model and its architecture we can improve the result further. Figure 7 shows the result of the modelling for test set.
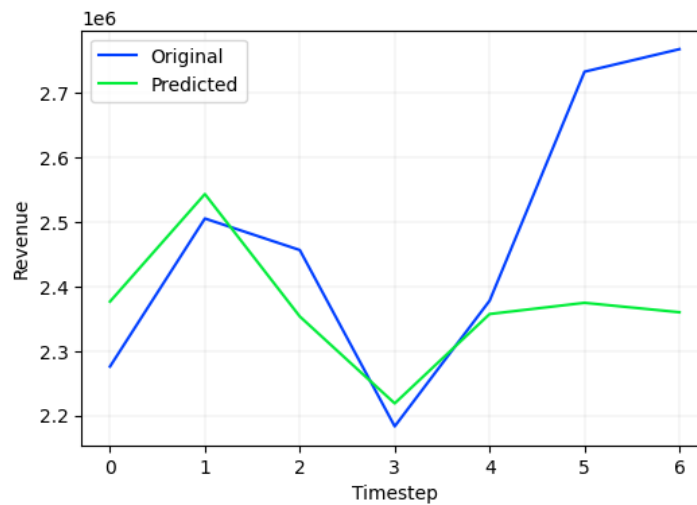


Figure 7: Deep Learning Results

# 2    Bayesian Methods

In this section I fit a simple neural network in a Bayesian framework using non-informative priors. In this example, the priors will be non-informative, which means they will have minimal impact on the final results and will let the data speak for itself. Non-informative priors are often used when we have little or no prior knowledge about the parameters we're trying to estimate, which is a good starting point in our Bayesian analysis.

Common noninformative priors include $unif(-1000, 1000)$ and $N(0, 10,000)$. In this task I used a normal distribution for the priors, specifically, for weight and biases of the model, I used a normal distribution with mean zero and standard deviation of 10.

A prior distribution is said to be conjugate to a likelihood function if the resulting posterior distribution is in the same family of probability distributions as the prior. For example, in the case of a normal likelihood function with unknown mean and known variance, the conjugate prior is another normal distribution. If the likelihood is a binomial distribution, the conjugate prior is a beta distribution. If the likelihood is a Poisson distribution, the conjugate prior is a gamma distribution. Looking at the posterior distribution I can say the priors are conjugate. For MCMC I used the default value in PyMC, which is NUTS, a sampler for continuous variables based on Hamiltonian mechanics. It allows us to sample from the posterior distribution and obtain estimates for the model parameters.

Comparing the results from the first section and the result of the Bayesian framework produced similar outputs. The model has the RMSE value of 94420.70 which is on par with the performance of the previous model used. Figure **??** shows the result of the prediction on the test set.
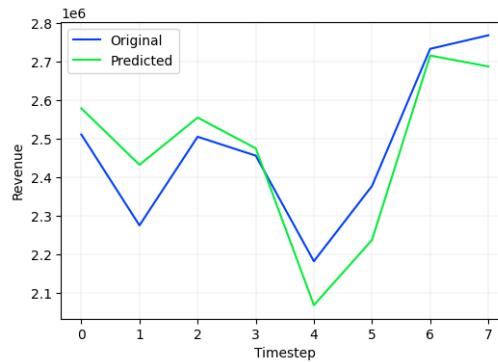


Figure 8: Bayesian Neural Network Result for Test Set

## 2.1    Prior choice

To modify the prior specification to account for the expert's belief that "Media 3 has no impact in generating marketing revenues," we can update the prior for the corresponding coefficient to reflect this belief. We want to incorporate the expert's opinion while still allowing the data to influence the final estimates. We can do this by using a prior distribution for the coefficient of "Media 3" with a high variance, indicating a wide range of possible values. For example, we can use a Normal distribution for the coefficient of Media 3 with a mean of 0 and a large standard deviation, such as 10. This would express the expert's belief that the coefficient is likely to be around 0, but they are not certain and allow the data to adjust this belief.

However, I already use a normal distribution with mean 0 and standard deviation of 10 to accommodate non informative priors in the model. So to answer the question raised in the task, I must say based on the estimated parameters from the model, the expert is correct and Media 3 has no to little impact on the revenue.

Additionally, to incorporate the opinions of the other expert, we can set informative priors for the coefficients of "newsletter," "media1_S," "media2_S," and "competitor_sales" based on the expert's beliefs. For example:
"newsletter": We can use a Normal distribution with a mean of 1 and a small standard deviation, such as 0.1, to express the belief that there is a strictly positive relationship between "newsletter" and marketing revenues. For "media1_S" we can use a Normal distribution with a mean of 2 and a small standard deviation, such as 0.1, to express the belief that one unit invested in "Media 1" will result in two units of marketing revenues returned. For "media2_S" we can use a Normal distribution with a mean of 8 (four times the mean for "media1_S") and a small standard deviation, such as 0.1, to express the belief that the impact of "Media

2" is four times that of "Media 1."

For "competitor_sales" we can use a Normal distribution with a mean of 0.15 and a small standard deviation, such as 0.05, to express the belief that for every unit increase in "competitor_sales," the change on revenues returned must range between 0 and 0.3.

Each of these hypothesis need to be implemented in the model separately and the results and posterior distribution needs to be carefully study. Unfortunately, I did not had enough time to do this analysis. The code needed for doing this is provided in the submitted code. Given more time, I could have make the necessary analysis to come up with answers.

# 3  Closing Remarks

I would want to convey my gratitude for the opportunity to demonstrate my qualifications for this post at the University of Essex and Croud Inc. This is an intriguing role for me, and I would be pleased to go to the next stage of the hiring procedure.

The Python codes can be found on the following pages. You will also receive a link to a colab notebook along with this report.
Click here to open in Google Colab

# Python Code

```python
# -*- coding: utf-8 -*-
"""alotfipoor_ktp_associate_essex.ipynb

Automatically generated by Colaboratory.

Original file is located at
    https://colab.research.google.com/drive/16cX08SFB-SyHqLJpUx33vgg6iqRHmVel

# Setup
"""

# Commented out IPython magic to ensure Python compatibility.
# from datetime import datetime
import os, glob, re
import numpy as np
import pandas as pd
import warnings
from math import sqrt

warnings.simplefilter(action='ignore', category=FutureWarning)
warnings.simplefilter(action='ignore', category=DeprecationWarning)

# import seaborn as sns
import matplotlib as mpl
import matplotlib.pyplot as plt
# %matplotlib inline
import seaborn as sns
from matplotlib.ticker import FuncFormatter
import matplotlib.dates as mdates

from IPython.core.display import display, HTML
display(HTML("<style>.container { width:100% !important; }</style>"))

# plotly
import plotly.express as px
import plotly.graph_objs as go
import plotly.io as pio
pio.renderers.default = 'colab'

from statsmodels.tsa.stattools import adfuller
from statsmodels.tsa.arima.model import ARIMA
from statsmodels.graphics.tsaplots import plot_pacf
from statsmodels.graphics.tsaplots import plot_acf
from statsmodels.tsa.seasonal import seasonal_decompose

import pymc as pm
import arviz as az

# import optuna

# from scipy.stats import uniform, randint
# from sklearn.model_selection import TimeSeriesSplit, cross_val_score, GridSearchCV, RandomizedSearchCV

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.linear_model import LinearRegression, SGDRegressor
from sklearn.svm import SVR
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import WhiteKernel, DotProduct
from sklearn.tree import DecisionTreeRegressor
```

```python
61  from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
62  from sklearn.neural_network import MLPRegressor
63  from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score, mean_squared_log_error
64
65  from keras.models import Sequential
66  from keras.layers import Dense, LSTM, Dropout, TimeDistributed, Flatten, BatchNormalization
67  from tensorflow.keras.optimizers import Adam
68  from tensorflow.keras.callbacks import LearningRateScheduler, EarlyStopping, ReduceLROnPlateau, ModelCheckpoint
69  # import tensorflow as tf
70
71  # data Table Display
72  # %load_ext google.colab.data_table
73
74  # to make this notebook's output stable across runs
75  np.random.seed(42)
76
77  # mount google colab
78  from google.colab import drive
79  drive.mount('/content/drive')
80
81  # load  dataset
82  df = pd.read_csv('/content/drive/MyDrive/[04] Colab Notebooks/data/weekly_media_sample.csv')
83  df['DATE'] = pd.to_datetime(df['DATE'])
84
85  df.head(5)
86
87  fig = go.Figure()
88  fig.add_trace(go.Scatter(x=df.DATE, y=df.revenue, name='Revenue', line=dict(color='blue', width=1)))
89  fig.add_trace(go.Scatter(x=df.DATE, y=df.competitor_sales, name='Competitor Sales', line=dict(color='red',
    ↪  width=.8)))
90  fig.add_trace(go.Scatter(x=df.DATE, y=df.newsletter, name='Newsletter Subscription', line=dict(color='green',
    ↪  width=.8)))
91  fig.show()
92
93  """# Part 1: Classical Marketing Data Modelling"""
94
95  # Creating a subset of dataframe for modelling
96  # df2 = df[['DATE','revenue','media1_S','media2_S','media3_S','competitor_sales','newsletter']].copy()
97  df2 = df[['DATE','revenue','media1_S','competitor_sales','newsletter']].copy()
98
99  # Calculate correlation coefficients with the target variable 'revenue'
100 correlation_matrix = df2.corr()
101 correlation_with_revenue = correlation_matrix['revenue'].abs().sort_values(ascending=False)
102
103 # Print the correlation coefficients
104 print(correlation_with_revenue)
105
106 mask = np.triu(np.ones_like(correlation_matrix, dtype=bool))
107 cmap = sns.diverging_palette(230, 20, as_cmap=True)
108 sns.heatmap(correlation_matrix,mask = mask, cmap=cmap, annot=True)
109
110 # Define the split date based on the given threshold (e.g., '2016-05-30' in this case)
111 split_date = '2016-05-30'
112
113 # Split the data into train and test sets
114 train = df2[df2.DATE <= split_date]
115 test = df2[df2.DATE > split_date]
116
117 # Extract the target variable 'revenue' from the training and test data
118 y_train = train['revenue']
119 y_test = test['revenue']
120
121 # Extract the input features from the training and test data
```

```
122  X_train = train.drop(['revenue', 'DATE'], axis=1)
123  X_test = test.drop(['revenue', 'DATE'], axis=1)
124
125  # Scale the data
126  scaler = StandardScaler()
127  X_train = scaler.fit_transform(X_train)
128  X_test = scaler.transform(X_test)
129
130  # Create instances of regression models with hyperparameters
131  linear_reg = LinearRegression()
132  sgd_reg = SGDRegressor()
133  # svm_reg = SVR(C=10, degree=3, gamma='auto', kernel='rbf')
134  # gaussian_process_reg = GaussianProcessRegressor(kernel=DotProduct() + WhiteKernel())
135  random_forest_reg = RandomForestRegressor(n_estimators=200, max_depth=20, random_state=42)
136  gradient_boosting_reg = GradientBoostingRegressor(n_estimators=200, learning_rate=0.01, max_depth=20,
     ↪  random_state=42)
137  mlp_reg = MLPRegressor(hidden_layer_sizes=(500, 300, 200), batch_size=64, learning_rate='adaptive', alpha=0.1,
     ↪  max_iter=5000, random_state=42)
138
139  # Create DataFrames to store predicted values
140  train_predictions_df = pd.DataFrame(data=y_train.values, columns=['Actual Revenue'])
141  test_predictions_df = pd.DataFrame(data=y_test.values, columns=['Actual Revenue'])
142
143  results_df = pd.DataFrame(columns=['Model', 'RMSE', 'MAE', 'R2', 'MAPE'])
144
145  # Mean Absolute Percentage Error (MAPE) function
146  def mean_absolute_percentage_error(y_true, y_pred):
147      return np.mean(np.abs((y_true - y_pred) / y_true)) * 100
148
149  # Function to evaluate and append results to the DataFrame
150  def evaluate_model(model_name, model, X_train, y_train, X_test, y_test):
151      model.fit(X_train, y_train)
152      y_train_pred = model.predict(X_train)
153      y_test_pred = model.predict(X_test)
154      train_predictions_df[model_name] = y_train_pred
155      test_predictions_df[model_name] = y_test_pred
156      rmse = mean_squared_error(y_test, y_test_pred, squared=False)
157      mae = mean_absolute_error(y_test, y_test_pred)
158      r2 = r2_score(y_test, y_test_pred)
159      mape = mean_absolute_percentage_error(y_test, y_test_pred)
160      results_df.loc[len(results_df)] = [model_name, rmse, mae, r2, mape]
161
162  # Evaluate each model and store the results
163  evaluate_model('Linear Regression', linear_reg, X_train, y_train, X_test, y_test)
164  evaluate_model('Stochastic Gradient Descent', sgd_reg, X_train, y_train, X_test, y_test)
165  # evaluate_model('SVM', svm_reg, X_train, y_train, X_test, y_test)
166  # evaluate_model('Gaussian Processes', gaussian_process_reg, X_train, y_train, X_test, y_test)
167  evaluate_model('Random Forest', random_forest_reg, X_train, y_train, X_test, y_test)
168  evaluate_model('Gradient Boosting', gradient_boosting_reg, X_train, y_train, X_test, y_test)
169  evaluate_model('Multi-layer Perceptron', mlp_reg, X_train, y_train, X_test, y_test)
170
171  results_df.head()
172
173  # Print the coefficient of the linear regression
174  linear_reg.coef_
175
176  # Merge train_predictions_df and test_predictions_df while keeping 'DATE' column
177  merged_predictions_df = pd.concat([train_predictions_df, test_predictions_df])
178
179  # Add the 'DATE' column from df2 to merged_predictions_df
180  merged_predictions_df['DATE'] = np.sort(df2['DATE'])
181
182  # Create a trace for each model's predicted values against the actual 'revenue'
```

```python
183  traces = []
184  for column in merged_predictions_df.columns[:-1]:  # Exclude the first and last columns ('Actual Revenue' and
     ↪    'DATE')
185      trace = go.Scatter(
186          x=merged_predictions_df['DATE'],
187          y=merged_predictions_df[column],
188          mode='lines',
189          name=column
190      )
191      traces.append(trace)
192
193  # Create the layout for the plot
194  layout = go.Layout(
195      title='Actual vs. Predicted Revenue',
196      xaxis=dict(title='Date'),
197      yaxis=dict(title='Revenue'),
198  )
199
200  # Create the figure and plot
201  fig = go.Figure(data=traces, layout=layout)
202  fig.show()
203
204  # Convert 'DATE' column to datetime object
205  # merged_predictions_df['DATE'] = pd.to_datetime(merged_predictions_df['DATE'])
206
207  # Create a line chart for each model's predicted values against the actual 'revenue' over time
208  plt.figure(figsize=(12, 6))
209  for column in merged_predictions_df.columns[1:-1]:  # Exclude the first and last columns ('Actual Revenue' and
     ↪    'DATE')
210      plt.plot(merged_predictions_df['DATE'], merged_predictions_df[column], label=column)
211
212  # Set the x-axis label to 'Date' and y-axis label to 'Revenue'
213  plt.xlabel('Date')
214  plt.ylabel('Revenue (Millions)')
215
216  # Set the title of the plot
217  plt.title('Actual vs. Predicted Revenue')
218
219  # Show the legend to identify each model's line
220  plt.legend()
221
222  # Format x-axis date labels to show only the month and year
223  date_format = mdates.DateFormatter("%b %Y")
224  plt.gca().xaxis.set_major_formatter(date_format)
225
226  # Set the x-axis tick interval to show labels every n months (adjust n as needed)
227  n_months = 2
228  plt.gca().xaxis.set_major_locator(mdates.MonthLocator(interval=n_months))
229
230  # Rotate the x-axis labels for better visibility
231  plt.xticks(rotation=45)
232
233  # Format y-axis tick labels in millions
234  def millions_formatter(x, pos):
235      return f'{x / 1e6:.1f}'
236
237  formatter = FuncFormatter(millions_formatter)
238  plt.gca().yaxis.set_major_formatter(formatter)
239
240  # Show the plot
241  plt.tight_layout()
242  plt.show()
243
```

```
244    # Fit the Linear Regression model
245    linear_reg.fit(X_train, y_train)
246
247    # Predict on the test set
248    Y_pred_linear = linear_reg.predict(X_test)
249
250    # Calculate the 95% confidence intervals for the predictions
251    conf_interval_linear = 1.96 * np.std(Y_pred_linear)  # 1.96 is the z-score for 95% confidence
252
253    # # Create a DataFrame with predictions and confidence intervals for the last 20 rows
254    # df_pred_linear = pd.DataFrame({'DATE': df2['DATE'].tail(8).values,
255    #                                'Actual Revenue': y_test[-8:].values.flatten(),
256    #                                'Predicted Revenue': Y_pred_linear[-8:].flatten(),
257    #                                'Lower CI': (Y_pred_linear[-8:] - conf_interval_linear).flatten(),
258    #                                'Upper CI': (Y_pred_linear[-8:] + conf_interval_linear).flatten()})
259
260    # Create a DataFrame with predictions and confidence intervals for the last 20 rows
261    df_pred_linear = pd.DataFrame({'DATE': df2['DATE'].tail(8).values,
262                                   'Actual Revenue': y_test.values.flatten(),
263                                   'Predicted Revenue': Y_pred_linear.flatten(),
264                                   'Lower CI': (Y_pred_linear - conf_interval_linear).flatten(),
265                                   'Upper CI': (Y_pred_linear + conf_interval_linear).flatten()})
266
267    # Plot the last 20 rows of predictions, actual observations, and confidence intervals
268    plt.figure(figsize=(12, 6))
269    plt.plot(df2['DATE'].tail(8), df2['revenue'].tail(8), label='Actual Revenue', color='blue')
270    plt.plot(df_pred_linear['DATE'], df_pred_linear['Predicted Revenue'], label='Predicted Revenue', color='green')
271    plt.fill_between(df_pred_linear['DATE'], df_pred_linear['Lower CI'], df_pred_linear['Upper CI'], alpha=0.2,
       ↪   color='orange')
272    plt.xlabel("Date")
273    plt.ylabel("Revenue")
274    plt.title("Linear Regression - Predicted vs. Actual (Test data) with 95% Confidence Intervals")
275    plt.legend()
276    plt.xticks(rotation=45)
277    plt.tight_layout()
278    plt.show()
279
280    """## ARIMA"""
281
282    sns_c = sns.color_palette(palette="deep")
283
284    sns.pairplot(
285        data=df2, kind="scatter", height=2, plot_kws={"color": sns_c[1]}, diag_kws={"color": sns_c[2]}
286    );
287
288    plt.figure(figsize=(12, 6))
289    plot_acf(df2['revenue'], lags=50)
290    plt.xlabel('Lag')
291    plt.ylabel('Autocorrelation')
292    plt.title('Autocorrelation Plot for Revenue')
293    plt.show()
294
295    plt.figure(figsize=(12, 6))
296    plot_pacf(df2['revenue'], lags=50)
297    plt.xlabel('Lag')
298    plt.ylabel('Partial Autocorrelation')
299    plt.title('Partial Autocorrelation Plot for Revenue')
300    plt.show()
301
302    # Performing seasonal decomposition
303    result = seasonal_decompose(df2['revenue'], model='additive', period=52)  # Assuming weekly seasonality
       ↪   (period=52)
304
```

```python
305    # Creating a subplot layout
306    fig, axes = plt.subplots(4, 1, figsize=(12, 10), sharex=True)
307
308    # Plotting the original time series
309    axes[0].plot(df2['DATE'], df2['revenue'], label='Original', color='blue')
310    axes[0].set_ylabel('Revenue')
311    axes[0].set_title('Original Time Series')
312
313    # Plotting the trend component
314    axes[1].plot(df2['DATE'], result.trend, label='Trend', color='orange')
315    axes[1].set_ylabel('Trend')
316    axes[1].set_title('Trend Component')
317
318    # Plotting the seasonal component
319    axes[2].plot(df2['DATE'], result.seasonal, label='Seasonal', color='green')
320    axes[2].set_ylabel('Seasonal')
321    axes[2].set_title('Seasonal Component')
322
323    # Plotting the residual component
324    axes[3].plot(df2['DATE'], result.resid, label='Residual', color='red')
325    axes[3].set_xlabel('Date')
326    axes[3].set_ylabel('Residual')
327    axes[3].set_title('Residual Component')
328
329    # Format x-axis date labels to show only the month and year
330    date_format = mdates.DateFormatter("%b %Y")
331    axes[-1].xaxis.set_major_formatter(date_format)
332
333    # Set the x-axis tick interval to show labels every n months (adjust n as needed)
334    n_months = 2
335    axes[-1].xaxis.set_major_locator(mdates.MonthLocator(interval=n_months))
336
337    # Rotate the x-axis labels for better visibility
338    plt.xticks(rotation=45)
339
340    # Adjusting the layout
341    plt.tight_layout()
342
343    # Show the plot
344    plt.show()
345
346    # Perform the ADF test on the 'revenue' data
347    result = adfuller(df2['revenue'])
348
349    # Extract the test statistic and p-value
350    test_statistic = result[0]
351    p_value = result[1]
352
353    print("ADF Test Statistic:", test_statistic)
354    print("p-value:", p_value)
355
356    if p_value <= 0.05:
357        print("The data is stationary.")
358    else:
359        print("The data is not stationary.")
360
361    # Find the optimal ARIMA hyperparameters using grid search
362    best_aic = float("inf")
363    best_order = None
364
365    # Define the ranges for p, d, and q
366    p_range = range(0, 15)  # Example: p can be 0, 1, or 2
367    d_range = range(0, 2)  # Example: d can be 0 or 1
```

```python
368    q_range = range(0, 3)   # Example: q can be 0, 1, or 2
369
370    for p in p_range:
371        for d in d_range:
372            for q in q_range:
373                try:
374                    arima_model = ARIMA(train['revenue'], order=(p, d, q))
375                    arima_fit = arima_model.fit()
376
377                    # Calculate AIC score for the current ARIMA model
378                    aic = arima_fit.aic
379
380                    if aic < best_aic:
381                        best_aic = aic
382                        best_order = (p, d, q)
383
384                except:
385                    continue
386
387    # Fit the ARIMA model with the best parameters to the training data
388    arima_model = ARIMA(train['revenue'], order=best_order)
389    arima_fit = arima_model.fit()
390
391    # Make predictions on the test set
392    y_pred_arima = arima_fit.forecast(steps=len(test))
393
394    # Evaluate the performance using Mean Squared Error (MSE)
395    mse_arima = mean_squared_error(test['revenue'], y_pred_arima)
396    rmse_arima = np.sqrt(mse_arima)
397
398    print("Best ARIMA Order (p, d, q):", best_order)
399    print("ARIMA RMSE:", rmse_arima)
400
401    # Plot the actual revenue and ARIMA predictions
402    plt.figure(figsize=(10, 6))
403    plt.plot(df2.index, df2['revenue'], label='Actual Revenue', color='blue')
404    plt.plot(test.index, y_pred_arima, label='ARIMA Predictions', color='red')
405    plt.xlabel('Date')
406    plt.ylabel('Revenue')
407    plt.title('Actual Revenue vs. ARIMA Predictions')
408    plt.legend()
409    plt.show()
410
411    # summary of fit model
412    print(arima_fit.summary())
413
414    # line plot of residuals
415    residuals = pd.DataFrame(arima_fit.resid)
416    residuals.plot()
417    plt.show()
418
419    # density plot of residuals
420    residuals.plot(kind='kde')
421    plt.show()
422    # summary stats of residuals
423    print(residuals.describe())
424
425    """## Deep learning"""
426
427    def create_time_series_dataset_with_lags(df, threshold_date, lag_observation=5, forecast_horizon=1):
428        """
429        Create X and Y matrices for time series modeling with lag observations and split into train and test sets.
430
```

```
431        Parameters:
432            df (pd.DataFrame): DataFrame containing the time series data with columns: revenue, competitor_sales,
433                               newsletter, Total_Media_Spend, and DATE.
434            threshold_date (str): The date to split the data into train and test sets (format: 'YYYY-MM-DD').
435            lag_observation (int): The number of lag observations to use. Default is 5.
436            forecast_horizon (int): The number of steps ahead to forecast revenue. Default is 1.
437
438        Returns:
439            X_train (np.ndarray): 3-dimensional array containing input features with lag observations for the train
    ↪   set.
440            Y_train (np.ndarray): 1-dimensional array containing target values (revenue) for the train set.
441            X_test (np.ndarray): 3-dimensional array containing input features with lag observations for the test
    ↪   set.
442            y_test (np.ndarray): 1-dimensional array containing target values (revenue) for the test set.
443        """
444        # Convert threshold_date to Timestamp object
445        threshold_date = pd.to_datetime(threshold_date)
446
447
448        # Get the 'revenue' time series and other features
449        revenue = df['revenue'].values
450        competitor_sales = df['competitor_sales'].values
451        newsletter = df['newsletter'].values
452        media1_S = df['media1_S'].values
453
454        X_train, Y_train, X_test, Y_test = [], [], [], []
455
456        for i in range(len(df) - lag_observation - forecast_horizon + 1):
457            # Input sequence with lag observations for all features
458            x_sequence = np.column_stack((
459                revenue[i:i + lag_observation],
460                competitor_sales[i:i + lag_observation],
461                newsletter[i:i + lag_observation],
462                media1_S[i:i + lag_observation]
463            ))
464
465            if df['DATE'][i + lag_observation - 1] <= threshold_date:
466                # Add to train set
467                X_train.append(x_sequence)
468                Y_train.append(revenue[i + lag_observation + forecast_horizon - 1])
469            else:
470                # Add to test set
471                X_test.append(x_sequence)
472                Y_test.append(revenue[i + lag_observation + forecast_horizon - 1])
473
474        X_train = np.array(X_train)
475        Y_train = np.array(Y_train)
476        X_test = np.array(X_test)
477        Y_test = np.array(Y_test)
478
479        return X_train, Y_train, X_test, Y_test
480
481    threshold_date = '2016-05-30'
482    lag_observation =3
483    forecast_horizon = 1
484
485    X_train, Y_train, X_test, Y_test = create_time_series_dataset_with_lags(df2, threshold_date, lag_observation,
    ↪   forecast_horizon)
486
487    # Normalize the input data
488    scaler = StandardScaler()
489    X_train_scaled = scaler.fit_transform(X_train.reshape(-1, X_train.shape[-1])).reshape(X_train.shape)
490    X_test_scaled = scaler.transform(X_test.reshape(-1, X_test.shape[-1])).reshape(X_test.shape)
```

```python
491
492    # Define the model
493    model = Sequential()
494    model.add(LSTM(128, activation='relu', return_sequences=True, input_shape=(X_train.shape[1], X_train.shape[2])))
495    model.add(LSTM(64, activation='relu', return_sequences=True))
496    model.add(Dropout(0.2))
497    model.add(Flatten())
498    model.add(Dense(128))
499    model.add(Dense(64))
500    model.add(Dense(64))
501    model.add(Dense(1))
502
503    model.compile(optimizer='adam', loss='mean_squared_error')
504
505    # Train the model
506    history = model.fit(X_train_scaled, Y_train, epochs=200, batch_size=16, verbose=0, shuffle=False)
507
508    # Make predictions
509    Y_pred = model.predict(X_test_scaled)
510
511    rmse = mean_squared_error(Y_test, Y_pred, squared = False)
512    print('RMSE for test dataset = %.5f' % rmse)
513
514    df_predicted = pd.DataFrame({'Actual Revenue': Y_test.flatten(), 'Predicted Revenue': Y_pred.flatten()})
515    print(df_predicted)
516
517    with plt.style.context('seaborn-bright', after_reset=True):
518      plt.figure(figsize=(6,4))
519      plt.plot(Y_test, label='Original')
520      plt.plot(Y_pred, label='Predicted')
521      plt.legend()
522      plt.xlabel('Timestep')
523      plt.ylabel('Revenue')
524      plt.grid(color='gray', linestyle='-', linewidth=0.1)
525      plt.show()
526
527    """# Part 2: Bayesian methods"""
528
529    # Define the split date based on the given threshold (e.g., '2016-05-30' in this case)
530    split_date = '2016-05-30'
531
532    # Split the data into train and test sets
533    train = df[df.DATE <= split_date]
534    test = df[df.DATE > split_date]
535
536    # Extract the target variable 'revenue' from the training and test data
537    y_train = train['revenue']
538    y_test = test['revenue']
539
540    # Extract the input features from the training and test data
541    x_train = train.drop(['revenue', 'DATE', 'X'], axis=1)
542    x_test = test.drop(['revenue', 'DATE', 'X'], axis=1)
543
544    # Scale the data
545    scaler = StandardScaler()
546    X_train = scaler.fit_transform(x_train)
547    X_test = scaler.transform(x_test)
548
549    # Define different priors for each feature's weights
550    priors = {
551        'media1_S': {
552            'mu': 0,
553            'sd': 10
```

```python
554          },
555          'media2_S': {
556              'mu': 0,
557              'sd': 10
558          },
559          'media3_S': {
560              'mu': 0,
561              'sd': 10
562          },
563          'competitor_sales': {
564              'mu': 0,
565              'sd': 10
566          },
567          'newsletter': {
568              'mu': 0,
569              'sd': 10
570          }
571      }
572
573      def contruct_bnn(x_input, y_input):
574          # Using the context manager to build the model
575          with pm.Model() as bnn:
576              # Data input
577              x_data = pm.MutableData('x_data', x_input)
578              y_data = pm.MutableData('y_data', y_input)
579
580              # Model structure
581              num_hidden = 5
582
583              # Priors for the weights
584              weights = {}
585
586              for feature in x_train.columns:
587                  weights[feature] = pm.Normal(f'w_{feature}', mu=priors[feature]['mu'], sigma=priors[feature]['sd'],
                     ↪  shape=num_hidden)
588
589
590              b_1 = pm.Normal('b_1', mu=0, sigma=10, shape=num_hidden)
591              w_out = pm.Normal('w_out', mu=0, sigma=10, shape=num_hidden)
592              b_out = pm.Normal('b_out', mu=0, sigma=10, shape=1)
593
594              # Neural network architecture
595              act_1 = pm.math.tanh(pm.math.dot(x_data, pm.math.stack([weights[feature] for feature in
                 ↪  x_train.columns])) + b_1)
596
597              act_out = pm.Deterministic('act_out', pm.math.dot(act_1, w_out) + b_out)
598
599              # Likelihood of the data (observed)
600              sigma = pm.HalfNormal('sigma', sigma=1)
601              likelihood = pm.Normal('Y_obs', mu=act_out, sigma=1, observed=y_data, total_size=y_input.shape[0])
602
603              return bnn
604
605      bnn = contruct_bnn(X_train, y_train)
606      # pm.model_to_graphviz(bnn)
607
608      # Specifying the MCMC algorithm (NUTS)
609      with bnn:
610          step = pm.NUTS()
611          trace = pm.sample(1000, tune=1000, step=step, chains=4, random_seed=42)
612
613      pm.plot_trace(trace, figsize=(14,10), legend=True, compact=False)
614      plt.show()
```

```
615
616   pm.set_data({"x_data": X_test, "y_data": y_test}, model=bnn)
617
618   # Generate posterior samples.
619   ppc_test = pm.sample_posterior_predictive(trace, model=bnn)
620
621   # Compute the point prediction by taking the mean
622   y_test_pred = ppc_test.posterior_predictive['Y_obs'].mean(dim=['chain', 'draw'])
623
624   # y_mean = y_test_pred.mean(axis=0)
625   rmse = mean_squared_error(y_test, y_test_pred, squared = False)
626   print('RMSE for test dataset = %.5f' % rmse)
627
628   with plt.style.context('seaborn-bright', after_reset=True):
629     plt.figure(figsize=(6,4))
630     plt.plot(y_test.values, label='Original')
631     plt.plot(y_test_pred.values, label='Predicted')
632     plt.legend()
633     plt.xlabel('Timestep')
634     plt.ylabel('Revenue')
635     plt.grid(color='gray', linestyle='-', linewidth=0.1)
636     plt.show()
637
638   pm.plot_trace(trace, figsize=(14,10), legend=True, compact=False)
639   plt.show()
640
641   pm.summary(trace, round_to=2)
642
643   # Plotting the ACF plot for the MCMC chains
644   pm.plot_autocorr(trace)
645   plt.show()
646
647   az.plot_posterior(trace)
648
649   """## Prior Choice"""
650
651   # Define different priors for each feature's weights
652   priors = {
653       'media1_S': {
654           'mu': 2,
655           'sd': 0.1
656       },
657       'media2_S': {
658           'mu': 8,
659           'sd': 0.1
660       },
661       'media3_S': {
662           'mu': 0,
663           'sd': 10
664       },
665       'competitor_sales': {
666           'mu': 0.15,
667           'sd': 0.05
668       },
669       'newsletter': {
670           'mu': 1,
671           'sd': 0.1
672       }
673   }
674
675   def contruct_bnn(x_input, y_input):
676       # Using the context manager to build the model
677       with pm.Model() as bnn:
```

```
678        # Data input
679        x_data = pm.MutableData('x_data', x_input)
680        y_data = pm.MutableData('y_data', y_input)
681
682        # Model structure
683        num_hidden = 1
684
685        # Priors for the weights
686        weights = {}
687
688        for feature in x_train.columns:
689            weights[feature] = pm.Normal(f'w_{feature}', mu=priors[feature]['mu'], sigma=priors[feature]['sd'],
               ↪  shape=num_hidden)
690
691
692        b_1 = pm.Normal('b_1', mu=0, sigma=10, shape=num_hidden)
693        w_out = pm.Normal('w_out', mu=0, sigma=10, shape=num_hidden)
694        b_out = pm.Normal('b_out', mu=0, sigma=10, shape=1)
695
696        # Neural network architecture
697        act_1 = pm.math.tanh(pm.math.dot(x_data, pm.math.stack([weights[feature] for feature in
               ↪  x_train.columns])) + b_1)
698
699        act_out = pm.Deterministic('act_out', pm.math.dot(act_1, w_out) + b_out)
700
701        # Likelihood of the data (observed)
702        sigma = pm.HalfNormal('sigma', sigma=1)
703        likelihood = pm.Normal('Y_obs', mu=act_out, sigma=1, observed=y_data, total_size=y_input.shape[0])
704
705        return bnn
706
707 bnn = contruct_bnn(X_train, y_train)
708 # pm.model_to_graphviz(bnn)
709
710 # Specifying the MCMC algorithm (NUTS)
711 with bnn:
712     step = pm.NUTS()
713     trace = pm.sample(1000, tune=1000, step=step, chains=4, random_seed=42)
714
715 pm.set_data({"x_data": X_test, "y_data": y_test}, model=bnn)
716
717 # Generate posterior samples.
718 ppc_test = pm.sample_posterior_predictive(trace, model=bnn)
719
720 # Compute the point prediction by taking the mean
721 y_test_pred = ppc_test.posterior_predictive['Y_obs'].mean(dim=['chain', 'draw'])
722
723 # y_mean = y_test_pred.mean(axis=0)
724 rmse = mean_squared_error(y_test, y_test_pred, squared = False)
725 print('RMSE for test dataset = %.5f' % rmse)
```