INTERNATIONAL
HELLENIC
UNIVERSITY

# Programmatic Automation and Yield Optimization on the Ad Exchange

## Galinos Giaglis

SID: 3305150010

Supervisor:                    Prof. Christos Tjortjis

SCHOOL OF SCIENCE & TECHNOLOGY

A thesis submitted for the degree of

*Master of Science (MSc) in E-Business, Innovation & Entrepreneurship (Part-time)*

DECEMBER 2018

THESSALONIKI – GREECE

# Abstract

This dissertation was written as a part of the MSc in eBusiness, Innovation & Entrepreneurship at the International Hellenic University.

The total media ad spending worldwide will rise to 7.4%, to $628 Bn in 2018, according to an eMarketer report, while the digital media (digital advertising) itself will account for almost 45% of the investments made, partly thanks to the global ecommerce shifting sizeable amounts of budgets from the TV ecosystem to the Digital placements. Marketer projections put digital advertising to a valuation of $129 Bn by 2021, making the digital advertising sector one of the largest in the non-tangible products market and one of the most intriguing to further investigate, automate and invest in. The advertising ecosystem is currently comprised of thousands of intermediary entities between an advertiser and a publisher; the two most essential aspects of this market. Unfortunately, the chain between the advertiser and the publisher is not a straight line and is usually filled with the intermediaries that in some cases provide excellent value and in other cases just intervene with the price and misalign the information between an advertiser and a publisher. The misinformation caused by the intermediaries along with the many different subchannels of digital advertising that exist in the market, e.g. direct buys, programmatic buys, performance buys, currently affect the way ads are being bought online. The purpose of this dissertation is to investigate the actions needed and create an algorithm that sets a (economic) yield management strategy, directly setting the price that they sell their ads on the global exchange.

This goal was successfully achieved by creating a software that works in four steps. First collects historical data from websites, their ad placements and the ad vendors' reporting tool, clears the redundant data, analyzes all possible variables affecting the price of the ad and deciding which truly affect the price and then using information around these variables makes suggestions for higher or lower price that can consecutively lead to higher revenue for a publisher.

I would like to acknowledge my supervisor Dr. Christos Tjortjis for his valuable help and support in all stages of this Dissertation.

Galinos Giaglis

7 December 2018

# Contents

# 1  Introduction

A publisher may decide to sell advertising space (ad inventory) online in two different and distinct ways. The first way it is by signing an Insertion Order (IO) with the Advertiser. In this IO the advertiser will buy a predefined amount of Ad Impressions (a term used for the quantity of advertising) at a predefined price in a predefined position on the website. This process is called "Direct Sell" or "Reservation Sell" and is considered the most sacred form of advertisement as each IO is as important as a contract between the advertiser and the publisher. When a publisher can't sell all of the ad inventory, one uses "Ad Exchanges", i.e Google, Appnexus, Rubicon Project, Index Exchange among others, which have several thousand buyers connected to them and act as an intermediary between their buyers and the publishers. By having such a large volume of advertisers, they have the ability to fill any gaps that exist from the publisher' inability to sell all of the ad inventory directly. This type of selling the **remnant inventory** is called "Programmatic Sell" and is the focus of investigation for this dissertation.

In our dissertation we will investigate the pricing strategies a publisher can set by utilizing a handful of variables that affect the price an ad can be sold at. Some of these variables are the size of the ad (e.g. 300x250, 300x600 etc.), the device type (e.g. Desktop, Mobile, Tablet, Connected TV), the transparency of the advertiser (e.g. Branded or Anonymous) among others. Our focus will be to identify all affecting variables, separate the most relevant, use the pertinent information to predict a higher revenue generating price and apply that strategy starting the cycle again the next day.

To be more specific, in this dissertation we present a time series algorithmic procedure on how to forecast floor pricing values of Google Ad Manager pricing rules, in order to achieve revenue increment (maximization) based on historical data. We describe our dataset extensively so that we can acquire information on the behavior of floor values along with the rest of the features and especially the ones highly correlated to the Ad eCPM. For each set of features that affect the Ad eCPM value, we tried several models in an iterative way in our effort to find the best fit to our dataset. These models were checked according to a set of statistical metrics and the results were visualized. Finally, for each set of features along with the model that fit best, we forecasted the values of Ad eCPM which indicated the optimal Floor value for each pricing rule. Furthermore, we are in the phase of developing a neural

network based on the proposed modeling approach that would programmatically learn, decide and set of the optimal Floor value of each pricing rule.

The first chapter of this dissertation includes the introduction where a brief walk through the algorithmic procedure that was developed is presented. In addition, a literature review is provided aside with the proposed approach definition. It focuses on the advantages that this approach has, including general information about the software developed

In the second chapter, an example dataset is analyzed and decomposed in order to understand the given dimensions and metrics that will try to fit in the modeling phase of the solution.

In the third chapter, the methodology of the analysis is explained including the five different models specified. In this section, each model of the analysis which concluded in using the Seasonal ARIMA process is explained with its weaknesses in comparison to the one produced by this process.

Next in the fourth chapter, the results and several outcomes of the analysis are shown. This section presents the results by visualizing the model fittings and the forecasting values aside with the original data. In addition, this section shows the final results after applying the software analysis and modeling to the real-time website pricing rules.

In the fifth chapter, the instructions of how to use the developed software are explained by giving a variety of ways to run it.

In the sixth chapter, the concluded results are shown with the suggestions for improvement.

The last chapter includes the full source code of the application in an ipython notebook format that ran on a specific website, over a time period of three months. The resulting revenue uplift is still under monitoring because of the exogenous factors that may apply to this optimization problem.

## 1.1  Literature overview

To the best of our knowledge, researchers' interest is focused on real time bidding (Wush Chi-Hsuan Wu, 2015), (Weinan Zhang, 2016), (Jun Wang, 2016) techniques which require data mining provided by scripting tools or custom pieces of code that reside on the client's website. These pieces of code are responsible for tracking features and metrics in order to understand and deliver statistical values and probabilistic models to solve the revenue optimization problem. Such solutions provide a state-of-the-art approach, in an academic interest perspective, but lack on application in real-life websites. Their major drawback rises when the owner of the website declines to provide access for equivalent research. Moreover, several solutions need to retain a waiting state of the ad unit in order to get the best revenue for each ad impression. This fact causes the "line-item to be auctioned" to expire, because of the time margin expiration leading to the impressions being lost.

Surveys like (Shuai Yuan, 2014) also state clearly how most of the website analysis are based on cookie manipulation and information retrieval in order to analyze efficiently online user characteristics so as to target respectively their ad campaigns.

## 1.2  Proposed approach definition

The proposed approach is trying to achieve a solution through experimentation with historical data retrieved from the ad vendor manager tool (Google Ad Manager) and not from the website directly. This gives us the opportunity to follow several dimensions and extract information based on a set of metrics which eventually provide the ability to forecast floor pricing values and thus optimize our revenue.

In addition, from a scientific point of view, we state that our algorithmic procedure tries to solve a time related problem, thus it requires a time series problem solving approach. This differentiates our solution from a regular regression problem solution in 2 ways:

1. It is ***time dependent***. So, the basic assumption of a linear regression model, that the observations are independent, doesn't hold in this case.

2. Along with an increasing or decreasing trend, most time series have some form of ***seasonality trends***, i.e. variations specific to a particular time frame. For example, if we visualize Ad Impressions over time, we will invariably find lower values in the weekdays rather than the weekend which depicts the fact that more people tend to browse their favorite content sites on the weekend and by that, the number of page requests increase which consequently increases the ad impressions.

Also, the resulted pricing rules are provided to the Google Ad Manager tool and processed for a period of three days over real-life websites in order to acquire feedback about the model parameters and the resulting revenue gain.
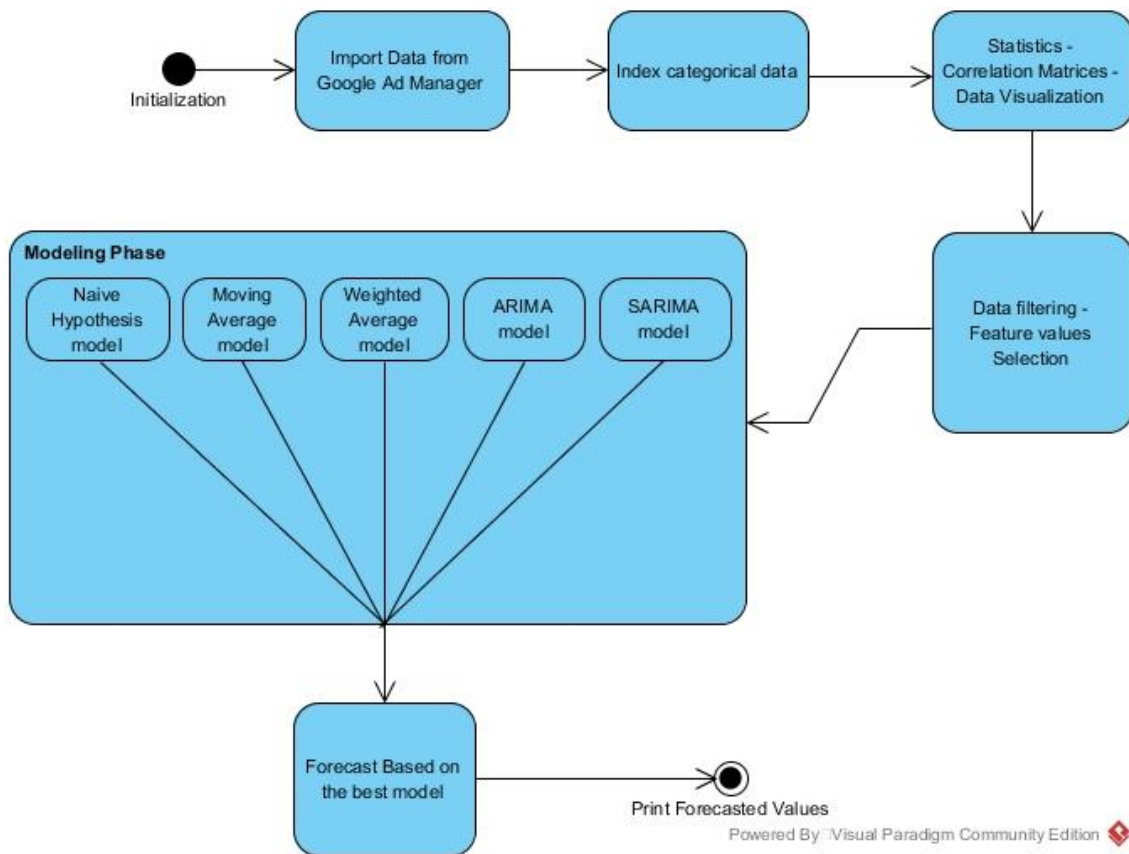
Figure 1: Flowchart diagram

Following the above assumptions, we propose an algorithmic procedure consisted of six steps as shown in Figure 1.

To start with, we acquire a historical query from the Google Ad Manager reporting tool containing specific dimensions and metrics as features of our dataset. These features will be manipulated and finally exposed to the modeling functions.

To derive numerical results from all the features gathered in the previous step, there has to follow an indexing step, where all the features that are not numerical will be transformed into numerical ones. Then the data set is ready to be studied through different statistical functions so that certain properties will be revealed. Among them are the correlation factors of the features such as inventory size, device category and branding type which will reveal association rules between them and will be used as filtering factors.

In the third step, the features that have the strongest correlation factors are chosen and specifically those that correlate best with Ad eCPM. Certain visualizations are provided to give the researcher a hint of the distribution of the data set values itself, along with histograms showing how close to the normal distribution these measures are. All features

[9]

compete at a correlation matrix having a score that indicates the correlation between their values, e.g., inventory size '*300x250*' is strongly correlated with mobile devices and especially with smart phones while their associated branding type is 'branded.'

While reaching step four, our data set is being filtered by the feature values that scored best and indicated from the third step. Each set of these pairs '*feature – value*' will be modeled independently and a different model will be trained by the software. The final set of feature values plus the labeled values is acquired for the modeling step aside with the aggregation information to form a per day value of Ad eCPM.

The fifth step is the modeling phase of this procedure and consists of six modeling solutions. Starting with a naïve hypothesis model, this procedure presents a baseline model. Next approaches are more sophisticated introducing a moving average model and as an extension of a weighted average model. A holt-winter exponential smoothing model is also demonstrated. At last an ARIMA model and a Seasonal ARIMA model present the best approach to this case study. Each of them are tested against several statistical measures and conclusively picked the best parameter initialization for them.

As a final step, this procedure uses the final model is selected on the modeling phaseand forecasts data for a period of time ahead. The resulting outcomes are visualized and demonstrated respectively along with evidence of correctness.

The proposed procedure by this dissertation unfortunately has also its drawbacks, with the most important to be introduced by the Google Ad Manager reporting tool which does not provide data with a sampling rate less than a day e.g. hourly or per minute. Thus, the presented procedure can forecast Ad eCPM values at a minimum per day time period basis. In case of a larger time period selected e.g. per week, a larger dataset should be acquired in order to limit the forecasting error inside the proposed probability confidence interval.

Moreover, exogenous factors that will influence Ad eCPM values in the 24 hours' time frame cannot be introduced into our model because of the lack of relevant information by the vendors' reporting tool.

## 1.3  Terms definition

The full meanings of the abbreviations used in this dissertation are:

Ad Manager – The Google tool used to deliver the ads on the page. It is essential an Ad Server or an intermediary delivering the ads from the publisher to the user

Ad Impressions – Absolute number that shows the amount of advertising ads shown

Ad Requests – Absolute number that shows the number an ad was requested (but not always shown)

Fill Rate or Coverage – The ratio of the Ad Impressions Served / the Ad Requests sent

Ad eCPM–The price for 1,000 Ad Impressions

CPM = Cost per Mile, the cost for 1,000 units

Revenue – Ad eCPM * (Ad Impressions/1000)

Ad Inventory – The available inventory the publisher can sell

Direct Sell – When a publisher sells his inventory directly to the advertiser

Programmatic Sell – When a publisher offers his remnant inventory in the Ad Exchanges

Branding Type – When the advertiser wishes to disclose her name and landing page URL. Two types of branding types exist, Branded (full disclosure) and Anonymous (no disclosure)

Floor Value – the minimum ad eCPM price required for an ad to be eligible for display in a website.

ACF – Autocorrelation function

PACF – Partial Autocorrelation function

AR – Autoregressive

MA – Moving Average

ARIMA – Autoregressive Integrated Moving Average

SARIMA – Seasonal Autoregressive Integrated Moving Average

AIC – Akaike Information Criterion

BIC – Bayesian Information Criterion

# 2 Data

## 2.1 Getting the Dataset

The required dataset was, at first, generated by using Google Ad Manager reporting tool where reports have been generated manually in Microsoft Excel sheets for each website that would be examined. After the finalization of this developer tool, the extraction of each dataset is going to be programmatically generated through the API of Google Ad Manager reporting tool.

Though the acquisition of the dataset came straight from the vendor's tool, the selection of the dimensions and measures that would successfully provide a well-formed and useful piece of data was rather difficult.

## 2.2 Describing the Dataset

Each row of the dataset consisted, as shown in Table 1, of four categorical fields (*Pricing rules, Inventory sizes, Device categories, and Branding types*) which represent the features, a date-time field which will be used later in the procedure as an index field and five numerical fields which represent the measures of each row.

Table 1: Data Sample before transformations

| | Days | Pricing rules | Inventory sizes | Device categories | Branding types | Ad requests | Ad impressions | Ad request eCPM (€) | Ad eCPM (€) | Diff AdCPM AdReqCPM | Estimated revenue (€) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2018-07-01 10:00:00 | GR Traffic | 300x250 | Connected TV | Anonymous | 326 | 290 | 0.094945 | 0.106731 | 0.011786 | 0.030952 |
| 1 | 2018-07-01 10:00:00 | GR Traffic | 300x250 | Connected TV | Branded | 100 | 87 | 0.371590 | 0.427115 | 0.055525 | 0.037159 |
| 2 | 2018-07-01 10:00:00 | GR Traffic | 300x250 | Desktop | Anonymous | 4046 | 3802 | 0.101428 | 0.107937 | 0.006509 | 0.410376 |
| 3 | 2018-07-01 10:00:00 | GR Traffic | 300x250 | Desktop | Branded | 1729 | 1630 | 0.350822 | 0.372130 | 0.021308 | 0.606571 |
| 4 | 2018-07-01 10:00:00 | GR Traffic | 300x250,320x100,320x50 | Connected TV | Anonymous | 4 | 4 | 0.106272 | 0.106272 | 0.000000 | 0.000425 |

Table 1 displays a raw sample of the data set provided from the reporting tool of Google Ad Manager before applying any transformation to the original data. It's worth mentioning that the information of the floor value of each pricing rule contained in the data set was obtained by a historical record of the website that was analyzed which was held in an external spreadsheet.

This difficulty was introduced by Google Ad Manager tool, which doesn't support the historical tracking of these values, but it can only provide the latest value used.

To support the proposed procedure, there has been an informal sub-step between the first step of importing data and the second step of indexing the categorical features, according to which the fields containing the information for date, pricing rule name and floor price value are matched between the two data frames.

In the indexing step of the procedure, the original data are transformed by an ordinal encoder transformation function, which is implemented in the Scikit–Learn toolkit (Scikit-Learn, 2018), to their numerical values respectively. This estimator transforms each categorical feature to one new feature of integer values starting from 0 to the number of distinct categories – 1. Such an integer representation can be used to convert categorical features to integer codes because their ordering is irrelevant to the information that they provide.

Table 2 shows a sample of the data set after the OrdinalEncoder transformation took place.

Table 2: Transformed sample of the data set through OrdinalEncoder

| Days | Ad requests | Ad impressions | Ad request eCPM (€) | Ad eCPM (€) | Floor | Diff AdCPM AdReqCPM | Estimated revenue (€) | Pricing rules | Inventory sizes | Device categories | Branding types |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2018-07-01 10:00:00 | 326 | 290 | 0.094945 | 0.106731 | 0.10 | 0.011786 | 0.030952 | 0.0 | 0.0 | 0.0 | 0.0 |
| 2018-07-01 10:00:00 | 100 | 87 | 0.371590 | 0.427115 | 0.35 | 0.055525 | 0.037159 | 0.0 | 0.0 | 0.0 | 1.0 |
| 2018-07-01 10:00:00 | 4046 | 3802 | 0.101428 | 0.107937 | 0.10 | 0.006509 | 0.410376 | 0.0 | 0.0 | 1.0 | 0.0 |
| 2018-07-01 10:00:00 | 1729 | 1630 | 0.350822 | 0.372130 | 0.35 | 0.021308 | 0.606571 | 0.0 | 0.0 | 1.0 | 1.0 |
| 2018-07-01 10:00:00 | 4 | 4 | 0.106272 | 0.106272 | 0.10 | 0.000000 | 0.000425 | 0.0 | 1.0 | 0.0 | 0.0 |

In order to get an insight of the given data set, the third step is dedicated to visualize the data distribution over all the features, categorical and numerical, against the Ad eCPM values.
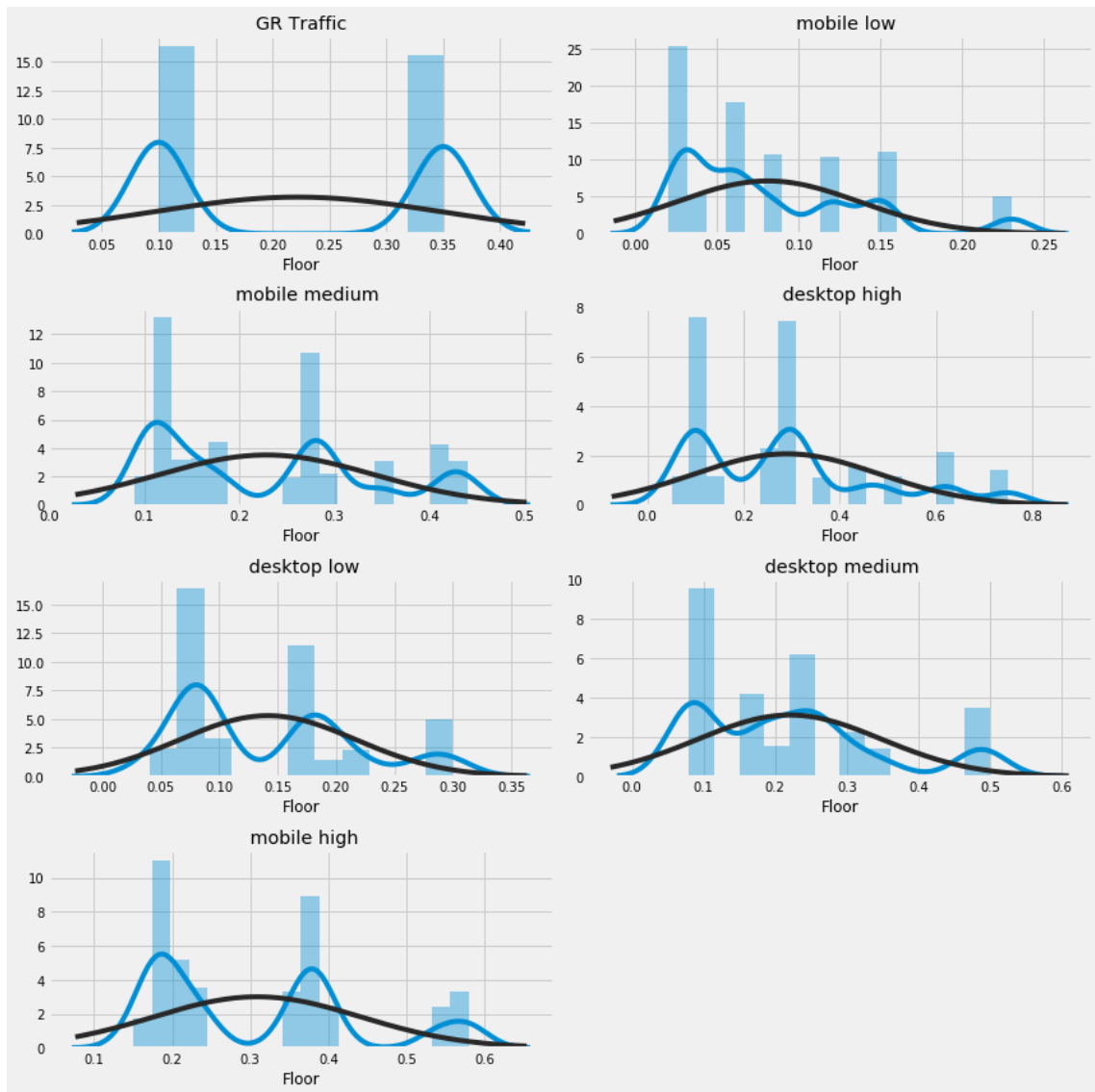
Figure 2: Floor pricing value distribution (blue line) with histogram
aside with normal distribution (black line)

Figure 2 represents the way that floor price values are distributed in each of the pricing rules specified for a certain website.

The following three figures give a detailed overview of how eCPM is distributed against the categorical features used in the data set.

Figure 3 shows that most of the generally well-paid ads are promoted by branding type '*Branded*' which is the most profitable of them regarding revenue.
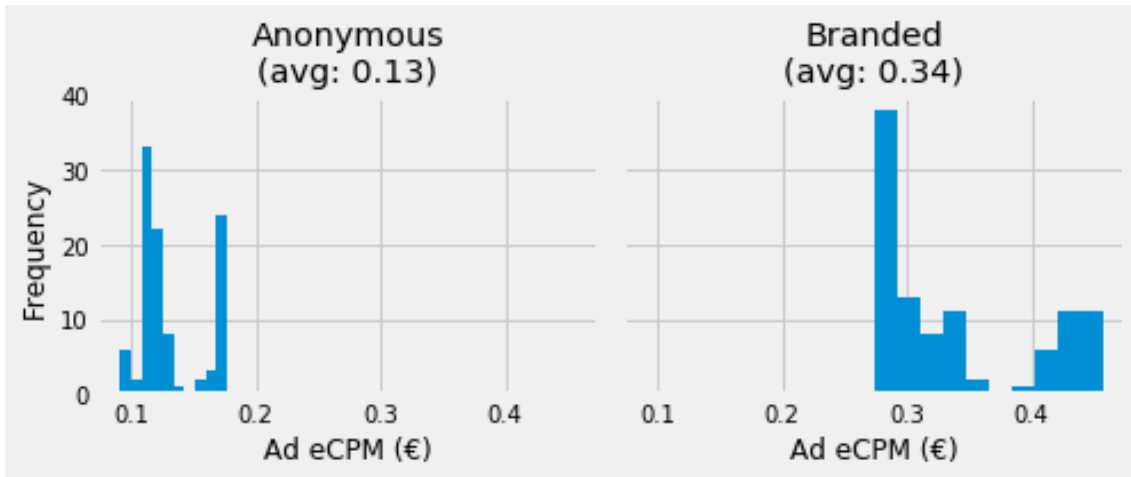
[15]

Figure 3: Branding Type histogram

Figure 4 represents the statistical analysis of the Device categories field of the data set, which shows that the best paid in average Ad eCPM are the ads that were displayed on a Desktop device while High-end mobile devices and Tablets share a rather large number of impressions.



Figure 4:Device Categories histogram

Moreover,

Figure 5 gives a detailed visual estimation of the importance of each inventory size and how it is being distributed to the different values of Ad eCPM. In the example presented, the most profitable size is '*300x250,320x100,320x50*' is presented with an average value of *0.28,* even though some of these sizes take part in other categories as well with lower average values.
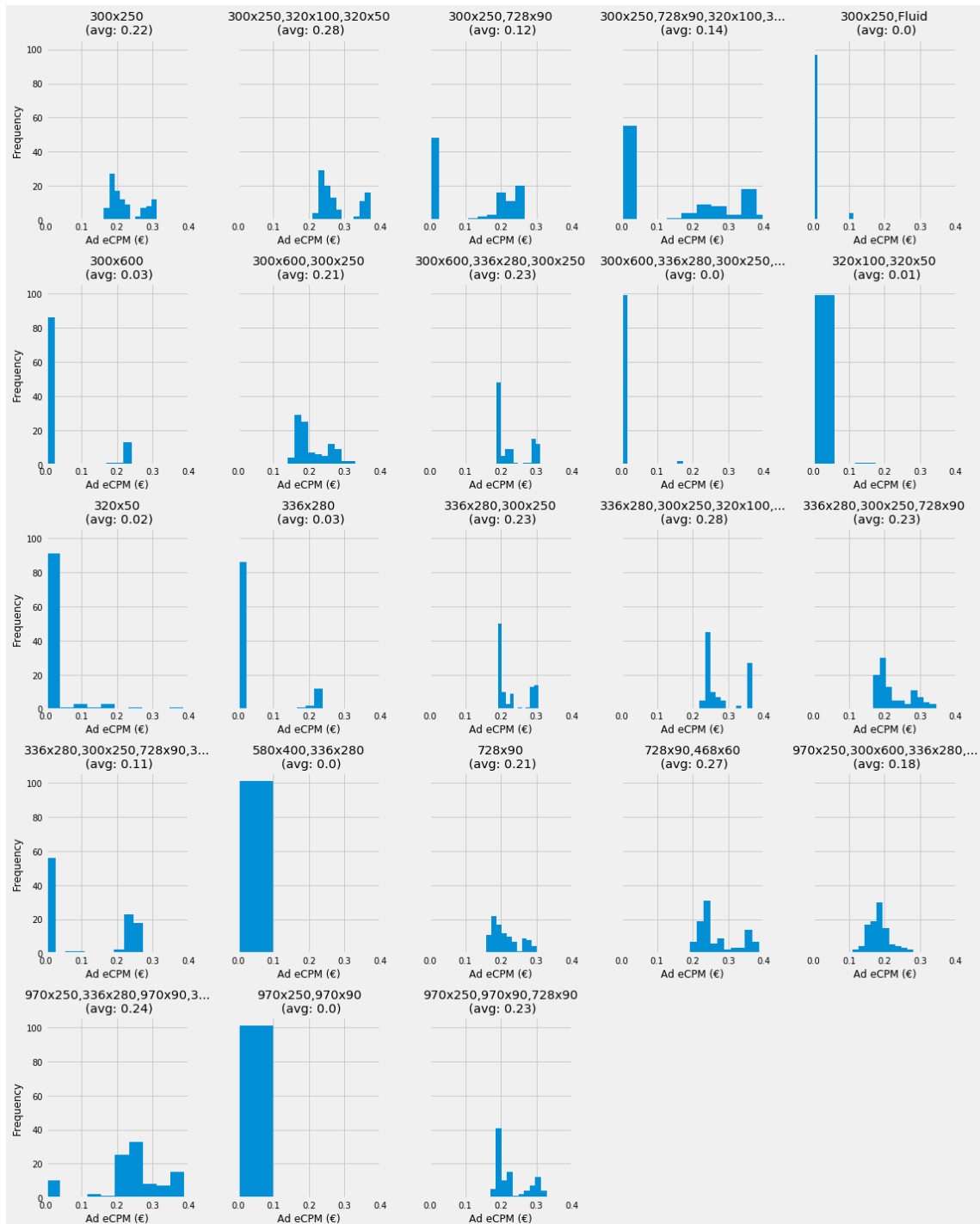


Figure 5: InventorySize Histogram

[17]

## 2.3 Categorical feature analysis

Supporting the observation that all the features are categorical, such as Pricing rules, Inventory sizes, Device Categories, and Branding types the proposed procedure indexes these features to their ordinal representation through the appropriate transformation and gives an integer value respectively. The rest data in the data set are numerical values that can participate in our modeling process as they are. Thus there is no need for extra preprocessing of the data.

In Figure 6 the categorical indexing is being displayed with the discrete values on the x-axis vs. the different features on the y-axis and the different pricing rules as the color coding.
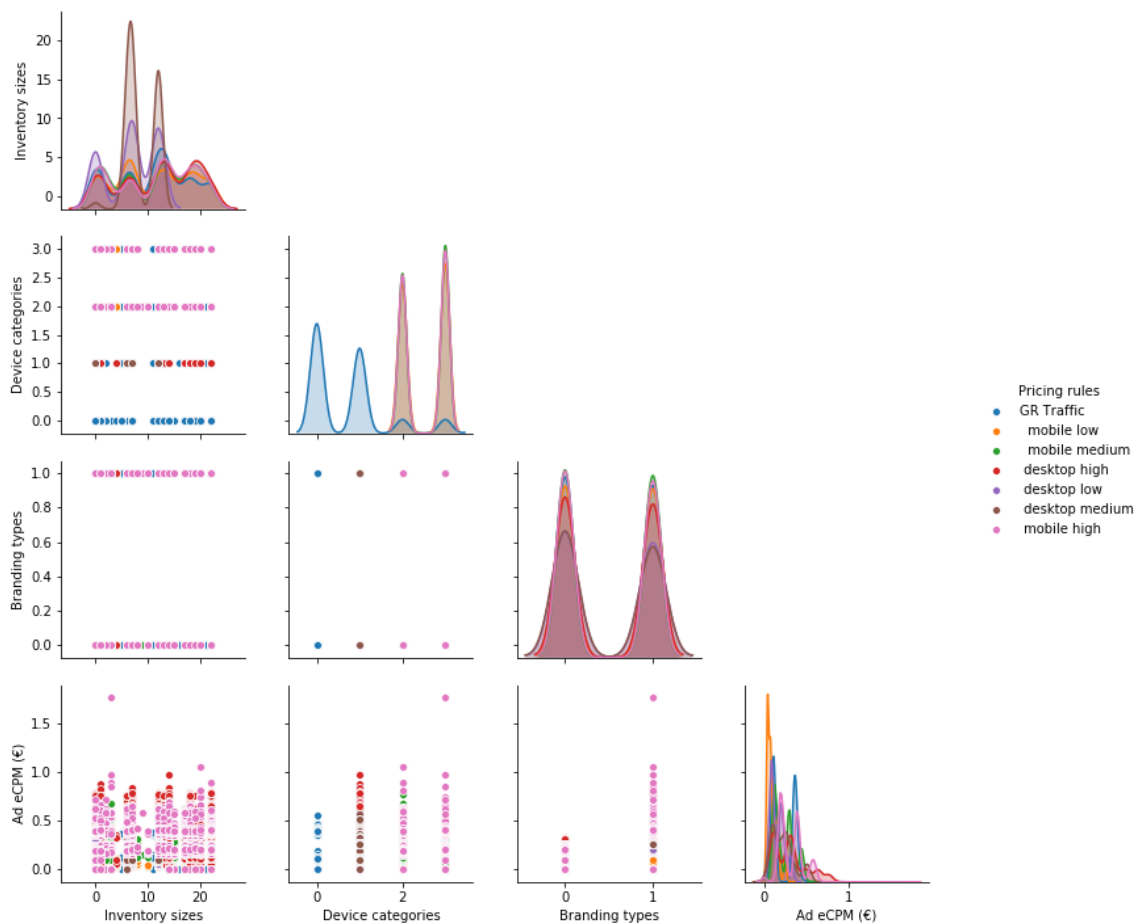


Figure 6: Categorical features distribution against Ad eCPM

[18]

## 2.4 Numerical measures analysis

In contrast with the categorical features, the numerical measures are imported as provided by the Google Ad Manager reporting tool.



Figure 7: Numerical features distribution
(*Ad eCPM, Ad request eCPM, Estimated revenue,* and *Floor values*)

Although it is mandatory to scale all continuous numeric input features so that not a single feature influences the model performance, we chose not to standardize our data set because the only features that could interfere with our modeling estimation are Ad requests and Ad Impressions.

Moreover, after the step of categorical features encoding, the data set visualization provides the capability to extract linearity between features to combine them in new ones properly. In this use case, the procedure showed a linear correlation between *Ad Requests – Revenue* and *Ad Impressions – Revenue* which can be explained because the Revenue value results from the multiplication of Ad Impressions and a value called *Close CPM* which depicts the actual payable

[19]

value of each impression. This multiplication gives an estimate of the revenue gained which on average is provided by Ad eCPM.

Also, the Ad Impressions measure can be described by a factor of Ad Request's value. This linearity between Ad requests, Ad Impressions, and Estimated revenue also denotes that these measures won't provide any information gain to the modeling procedure and their values will affect the revenue directly and not the Ad eCPM value.

As we can see both in Figure 7 and Figure 8, there is a strong statistical relationship which is measured by Pearson correlation factor (0.943 or 94.3%) between Ad eCPM and Floor value of each pricing rule which is the key property of our research.

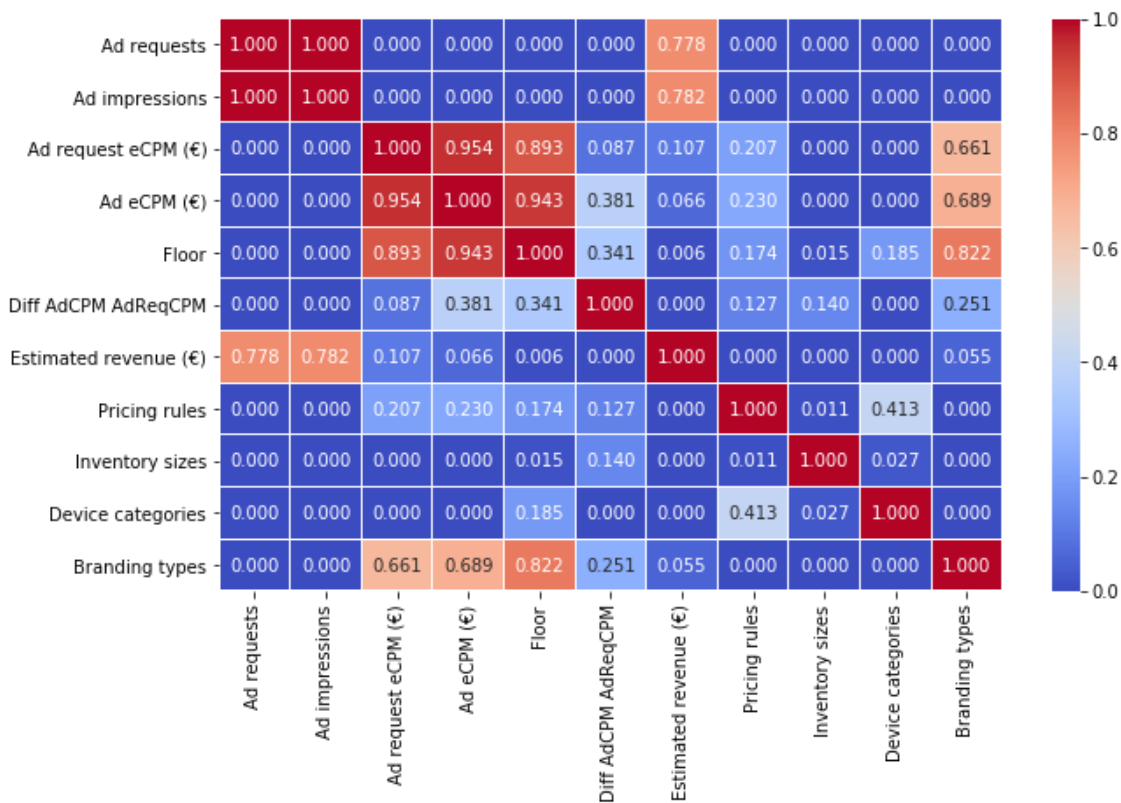| | Ad requests | Ad impressions | Ad request eCPM (€) | Ad eCPM (€) | Floor | Diff AdCPM AdReqCPM | Estimated revenue (€) | Pricing rules | Inventory sizes | Device categories | Branding types |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Ad requests | 1.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.778 | 0.000 | 0.000 | 0.000 | 0.000 |
| Ad impressions | 1.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.782 | 0.000 | 0.000 | 0.000 | 0.000 |
| Ad request eCPM (€) | 0.000 | 0.000 | 1.000 | 0.954 | 0.893 | 0.087 | 0.107 | 0.207 | 0.000 | 0.000 | 0.661 |
| Ad eCPM (€) | 0.000 | 0.000 | 0.954 | 1.000 | 0.943 | 0.381 | 0.066 | 0.230 | 0.000 | 0.000 | 0.689 |
| Floor | 0.000 | 0.000 | 0.893 | 0.943 | 1.000 | 0.341 | 0.006 | 0.174 | 0.015 | 0.185 | 0.822 |
| Diff AdCPM AdReqCPM | 0.000 | 0.000 | 0.087 | 0.381 | 0.341 | 1.000 | 0.000 | 0.127 | 0.140 | 0.000 | 0.251 |
| Estimated revenue (€) | 0.778 | 0.782 | 0.107 | 0.066 | 0.006 | 0.000 | 1.000 | 0.000 | 0.000 | 0.000 | 0.055 |
| Pricing rules | 0.000 | 0.000 | 0.207 | 0.230 | 0.174 | 0.127 | 0.000 | 1.000 | 0.011 | 0.413 | 0.000 |
| Inventory sizes | 0.000 | 0.000 | 0.000 | 0.000 | 0.015 | 0.140 | 0.000 | 0.011 | 1.000 | 0.027 | 0.000 |
| Device categories | 0.000 | 0.000 | 0.000 | 0.000 | 0.185 | 0.000 | 0.000 | 0.413 | 0.027 | 1.000 | 0.000 |
| Branding types | 0.000 | 0.000 | 0.661 | 0.689 | 0.822 | 0.251 | 0.055 | 0.000 | 0.000 | 0.000 | 1.000 |

Figure 8: Correlation Diagram against all dimensions

Based on this working hypothesis this algorithm provides a forecasting model of the Ad eCPM against other dimensions introduced by our data set, to achieve increased revenue with the appropriate adjustment of the Floor value of each pricing rule on a per-site basis.

# 3  Methodology

## 3.1  Description

As a starting step on the modeling phase of this algorithmic procedure, a correlation matrix – diagram as shown in Figure 9 is created, in order to model the Ad eCPM behavior. This matrix consists of Pearson correlation factor values against the information gained about Ad eCPM values.

Subsequently, feature sets that have strong correlation factors over the three major categories: Branding types – Device Categories – Inventory Sizes are defined. Each set consists of a tuple of three integer values, e.g. (1, 1, 7) which indicates that there is a strong relationship among Branded type equal to 1 – '*Branded*', Device category equals to 1 – '*Desktop*' and Inventory size equals 7 – '*300x600, 336x280, 300x250*'.
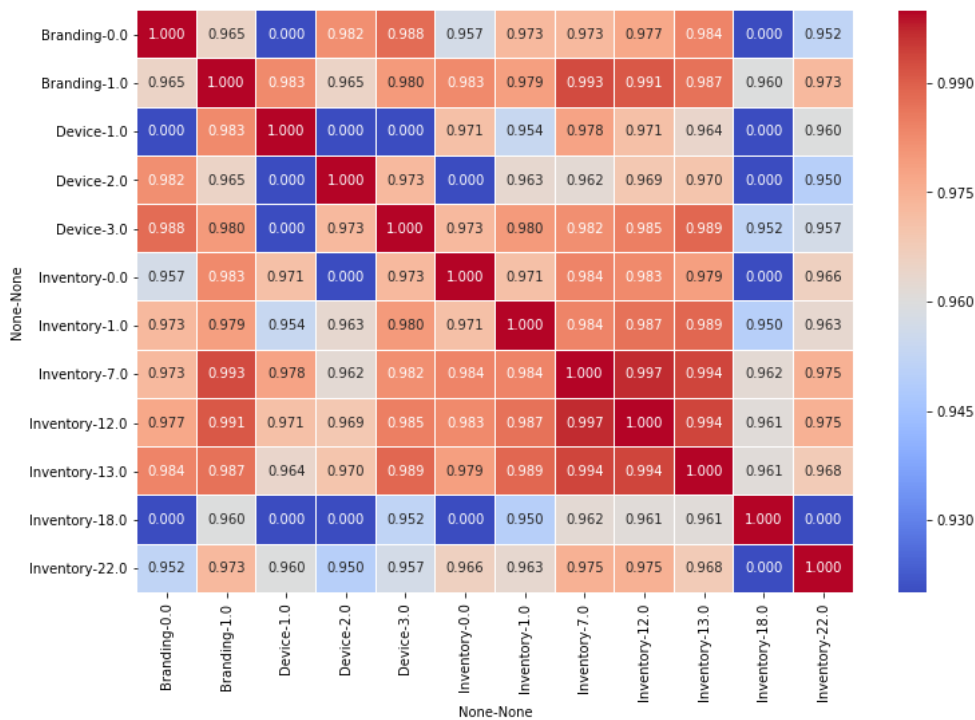


Figure 9: Correlation diagram based on average Ad eCPM values against
all dimensions of (Branding type, Device Category, Inventory Size)

Each of these sets underwent to a testing phase of different modeling schemes so that the best model fitting would be acquired.

[21]

## 3.2  Model Specification

To achieve the best fitting for the proposed model to the specified data set, five different modeling approaches were used, starting from a naïve hypothesis as our baseline model to more sophisticated models created by ARIMA and Seasonal ARIMA processes.

### 3.2.1  Naive hypothesis – Simple Average

A model like $\hat{y}_t = y_{t-1}$ is a great baseline for any time series prediction, as it relies to just the previous value of our data series and thus the error that will be introduced is rather small. Moreover, financial related series are likely to depend on the previous time period value and as forecasting steps increase, the error margin increase too. To overcome this problem, we, assumed that the future value of our variable depends on the average of its $k$ previous values and $k$ equals to the length of the time series. Such forecasting technique which forecasts the expected value equal to the average of all previously observed points is called Simple Average technique.

$$\hat{y}_t = \frac{1}{k} \sum_{n=1}^{k} y_{t-n}$$

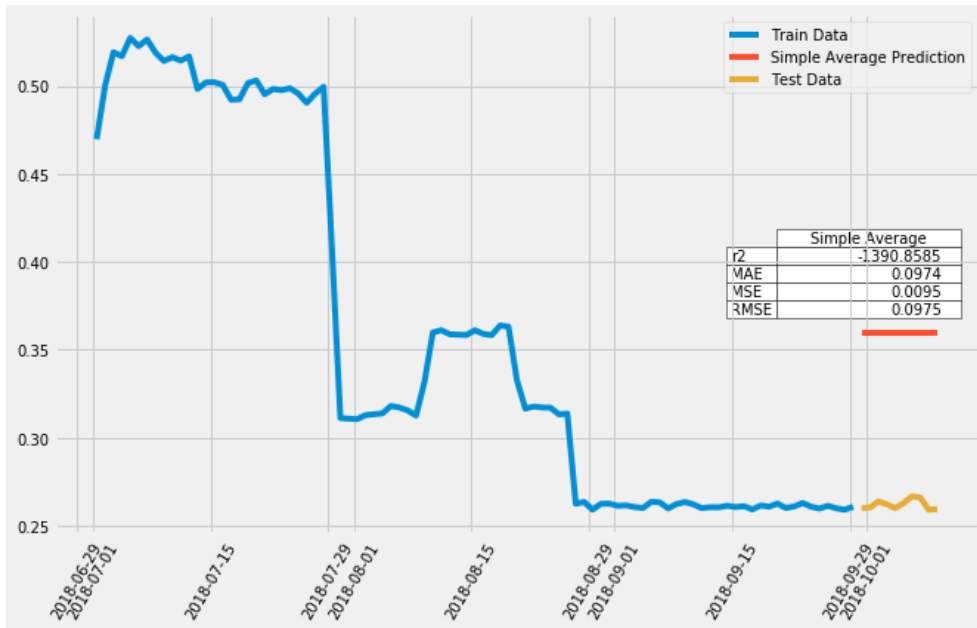| | Simple Average |
|---|---|
| r2 | -1,390.8585 |
| MAE | 0.0974 |
| MSE | 0.0095 |
| RMSE | 0.0975 |

Figure 10: Simple average modeling

Although this approach gives a prediction with a reasonable error for the first forecasting step, most of the time, if it is used to forecast more than one step ahead it will lead to a rather large amount of errors in the last forecasting value.

### 3.2.2 Moving average

To provide a more accurate forecasting value, we used a simple moving average model, which forecasts the next value(s) in a time series based on the average of a fixed finite number $p$ of the previous values. Thus, for all $i > p$

$$\hat{y}_t = \frac{1}{p} \sum_{n=1}^{p} y_{t-n}$$

This approach also incorporates the seasonal flow of the values of the time series, but again it only gives a rather rough prediction when we increment the forecasting steps.



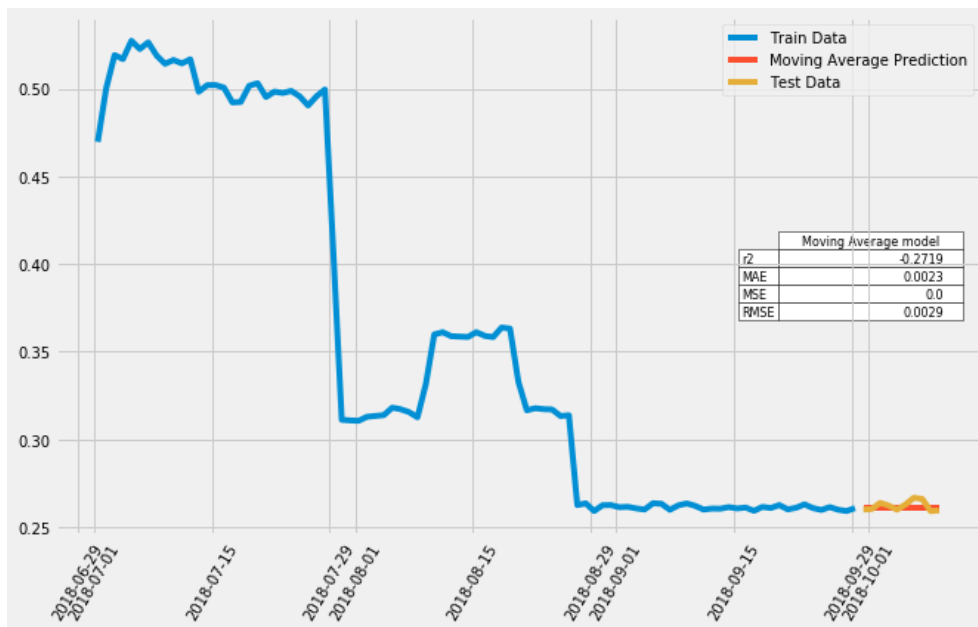| Moving Average model | |
| --- | --- |
| r2 | -0.2719 |
| MAE | 0.0023 |
| MSE | 0.0 |
| RMSE | 0.0029 |

Figure 11: Moving average example modeling

As shown in the above Figure 11 the moving average modeling method fits better the training dataset because it incorporates the seasonal trend of the data. The basic weakness of this modeling approach is that if the train data set has few data points; the error produced increments rapidly as forecasting steps increase.

### 3.2.3 Weighted moving average

As a next modeling approach, we chose a weighted moving average model which uses for the forecasting step value a different way of weighting the past observations and sum up to the weight value to one. The larger weights often assigned to the more recent observations promoting the corresponding values.

$$\hat{y}_t = \sum_{n=1}^{k} w_n y_{t+1-n}$$

In the above equation $w_n$ are the weights for each of the previous values and $k$ is the number of them to consider in the sum.



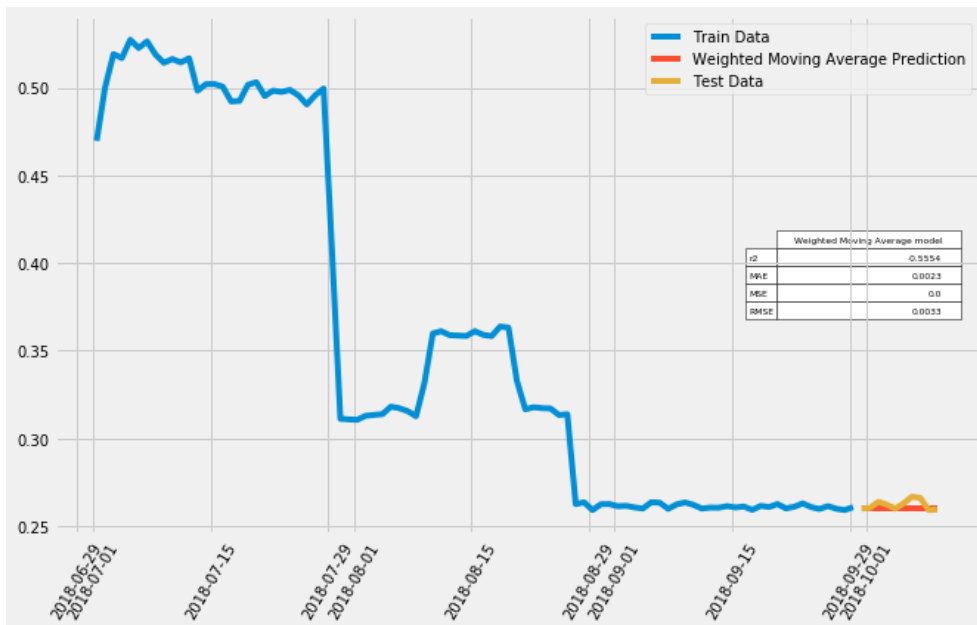| Weighted Moving Average model | |
|---|---|
| r2 | -0.5554 |
| MAE | 0.0023 |
| MSE | 0.0 |
| RMSE | 0.0033 |

Figure 12: Weighted moving average example modeling

This modeling approach has the advantage that if the autocorrelation factor of the data set to be fit depends over the k lag values, the predicted values approximate better the testing data set. Given this attribute, this approach also has its weakness because the weights array has a finite number of elements.

### 3.2.4 ARIMA (p, d, q)

Autoregressive Integrated Moving Average (ARIMA) models consist of 3 basic steps. The Auto regression part AR ($p$) where $p$ is the order of the AR model, the Integration I ($d$) part and the Moving average MA ($q$) part. While exponential smoothing models were based on a description of trend and seasonality in the data, ARIMA models aim to describe the correlations in the data with each other.

$$\Phi(B)(1 - B)^d y_t = c + \theta(B)\varepsilon_t$$

Where $\varepsilon_t$ is a white noise process with mean zero and variance $\sigma^2$, $B$ is the backshift operator, and $\Phi(z)$ and $\theta(z)$ are polynomials of order $p$ and $q$ respectively. If $c \neq 0$ there is an implied polynomial of order $d$ in the forecast function.

Figure 13 shows exactly the model predictions over the train data and the accuracy of each ARIMA model of order $p$ and $q$ respectively.
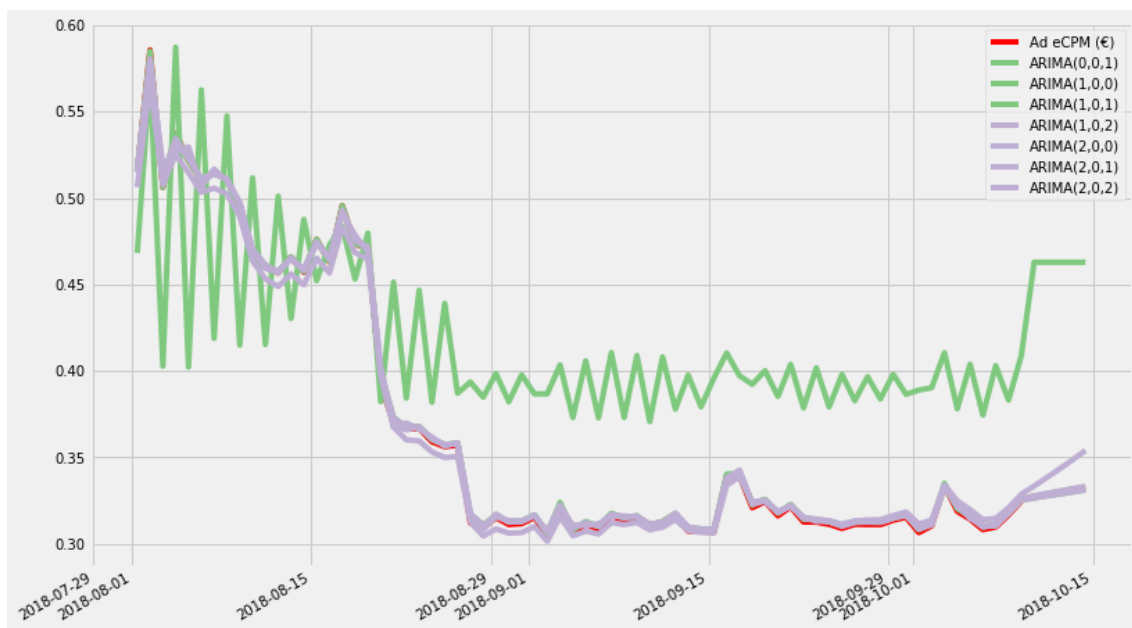


Figure 13: Different parameter values in ARIMA (p, d, q) along with the original data
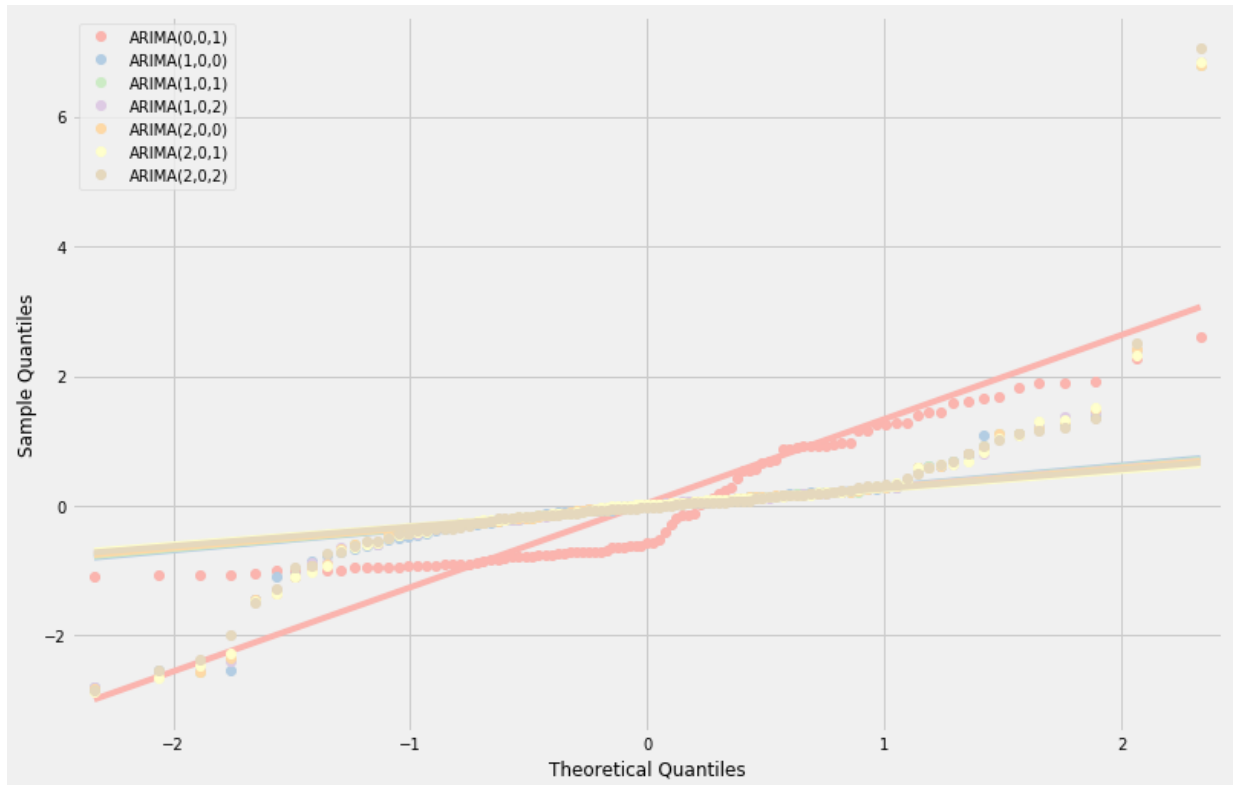
Figure 14: Quantiles convergence plots

In addition, in Figure 14 we present the plots that provide us information about the approximated values for the ARIMA modeling approach as described by *Algorithm 1*.

All of the above modeling schemes are all examined through statistical measures of error and information retrieval. One of them is the Mean Square Error rate (MSE), and another is the $r^2$ measure which provides the metric of the accuracy for the examined model.

The Akaike information criterion (AIC) (Bozdogan, 1987) is an estimator of the relative quality of statistical models for a given set of data. Given a collection of models for the data, AIC estimates the quality of each model, relative to each of the other models. Thus, AIC provides a means for model selection. The advantages are that it is valid for both nested and non-nested models, it can compare models with different error distribution and finally, it can avoid multiple testing issues.

Some weakness of the AIC is that it cannot be used to compare models of different data sets. Thus, the selected model with the lowest AIC is only valid and better than another for the specific data set.
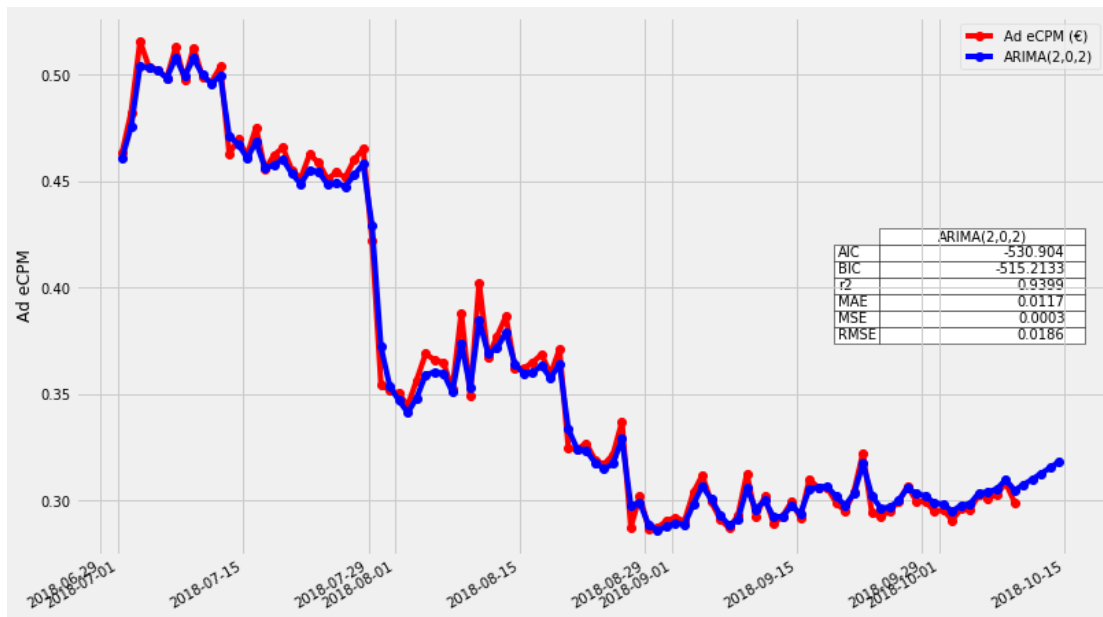
[27]

Figure 15 ARIMA Model fitting

As an example of the ARIMA model fitting Figure 15 show the resulting visualization for the predicted data denoted by the blue line against the original data set. As an extension of the information that describes best the model Table 3 presents thoroughly all the parameters required to define this model fitting.

Table 3: ARIMA (2, 0, 2) model summary

```
===============================================================================
                        coef    std err         z      P>|z|     [0.025     0.975]
-------------------------------------------------------------------------------
const                 0.3418      0.017    19.827      0.000      0.308      0.376
ar.L1.Ad eCPM (€)     1.9554      0.024    80.301      0.000      1.908      2.003
ar.L2.Ad eCPM (€)    -0.9576      0.025   -38.471      0.000     -1.006     -0.909
ma.L1.Ad eCPM (€)    -0.6248      0.096     -6.480      0.000     -0.814     -0.436
ma.L2.Ad eCPM (€)    -0.3752      0.093     -4.014      0.000     -0.558     -0.192
-------------------------------------------------------------------------------
```

[28]

### 3.2.5 Seasonal ARIMA (SARIMA)

An improvement over ARIMA is Seasonal ARIMA which takes into account the seasonality of dataset. As shown in Figure 16 the original data are decomposed into a series of trend and seasonality. This decomposition is required to determine the seasonality factor which will be used by the equation of the appropriate model.
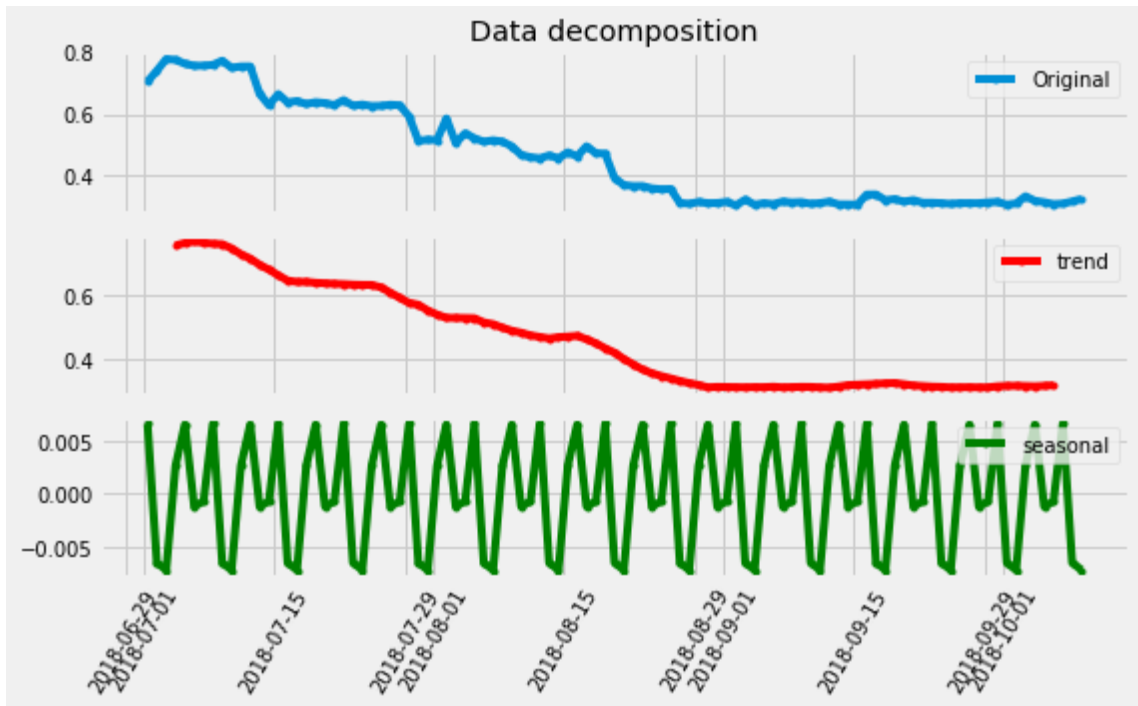


Figure 16: Data decomposition showing trend and seasonality

First of all, this class of models was introduced by Box and Jenkins (1976) and offers a good representation of many seasonal series that we find in practice and in simplified form is written as the ARIMA model $(P, D, Q)_m \times (p, d, q)$ where the period of the seasonal series is $m$.

$$\Phi(B^m)\varphi(B)(1 - B^m)^D(1 - B)^d y_t = c + \Theta(B^m)\theta(B)\varepsilon_t$$

$\Phi(z)$ and $\Theta(z)$ are polynomials, each containing no roots inside the unit circle, representing the seasonal AR operator of order $P$ and the seasonal moving average MA operator of order $Q$ respectively. If $c \neq 0$ the implied polynomial is of order $d + D$ in the forecasting function.

[29]

ARIMA and SARIMA forecasting is selecting an appropriate model order, that is the values $p, q, P, Q, D, d$. If $d$ and $D$ are known, we can select the orders $p, q, P$ and $Q$ via an information criterion such as the $AIC$:

$$AIC = -2 \log(L) + 2(p + q + P + Q + k)$$

Where $k = 1$ if $c \neq 0$ and $0$ otherwise, and $L$ is the maximized likelihood of the model fitted to the differenced data $(1 - B^m)^D (1 - B)^d y_t$ .

As already mentioned, the goal again is to select the model that minimizes the AIC amongst all of the models that are appropriate for the data.

The AIC (K.P. Burnham, 2004) also provides a method for selecting between the additive and multiplicative error models. The point forecasts from the two models are identical so that standard forecast accuracy measures such as the MSE or mean absolute percentage error (MAPE) are unable to select between the error types. The AIC is able to select between the error types because it is based on likelihood rather than one-step forecasts.
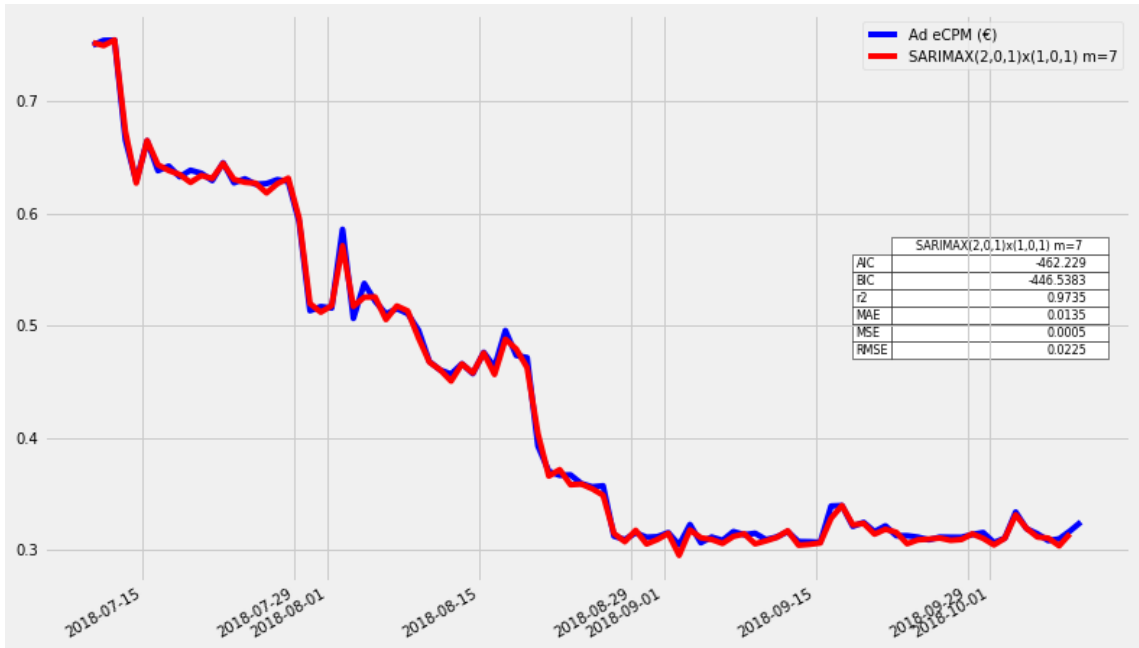


Figure 17 SARIMA model fitting

Again as an example of the SARIMA process of model fitting, Figure 17 visualizes the way that predictions are fit over the original data set, while

[30]

Table 4 describe the coefficients of the model function.

Table 4: Model coefficients description table

```
==============================================================================
                     coef    std err          z      P>|z|      [0.025      0.975]
------------------------------------------------------------------------------
const              0.4853      0.045     10.861      0.000       0.398       0.573
ar.L1.Ad eCPM (€)  1.9915      0.007    305.606      0.000       1.979       2.004
ar.L2.Ad eCPM (€) -0.9928      0.007   -151.739      0.000      -1.006      -0.980
ma.L1.Ad eCPM (€) -1.1556      0.113    -10.187      0.000      -1.378      -0.933
ma.L2.Ad eCPM (€)  0.1574      0.110      1.432      0.155      -0.058       0.373
==============================================================================
```

As an example of SARIMA$(2,0,1)_7(1,0,1)$ the above equations evaluate to:

$$(1 - B^7)^0(1 - B)^0 = 1$$

$$\Phi(B^m) = \Phi(B^7) = 1 - \Phi_1 B^7$$

$$\varphi(B) = 1 - \varphi_1 B$$

$$\Theta(B^m) = \Theta(B^7) = 1$$

$$\theta(B) = 1 + \theta_1 B$$

Forming the final equation of the model which is:

$$\Phi(B^m)\varphi(B)y_t = c + \theta(B)\varepsilon_t$$

$$(1 - \Phi_1 B^7)(1 - \varphi_1 B)y_t = c + (1 + \theta_1 B)\varepsilon_t$$

$$(1 - \varphi_1 B - \Phi_1 B^7 + \Phi_1 \varphi_1 B^8)y_t = c + (1 + \theta_1 B)\varepsilon_t$$

$$y_t - \varphi_1 y_{t-1} - \Phi_1 y_{t-7} + \Phi_1 \varphi_1 y_{t-8} = c + \varepsilon_t + \theta_1 \varepsilon_{t-1}$$

$$y_t = \varphi_1 y_{t-1} + \Phi_1 y_{t-7} - \Phi_1 \varphi_1 y_{t-8} + c + \varepsilon_t + \theta_1 \varepsilon_{t-1}$$

And finally, the forecasting function will be:

$$y_{t+1} = \varphi_1 y_t + \Phi_1 y_{t-6} - \Phi_1 \varphi_1 y_{t-7} + c + \theta_1 \varepsilon_t$$
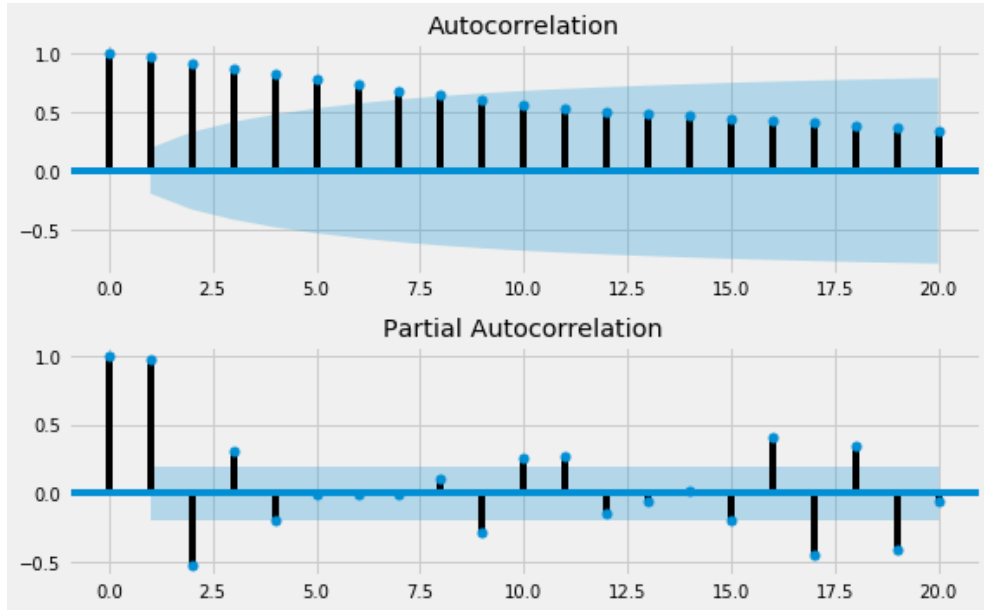
Figure 18: Autocorrelation and partial autocorrelation plotting of the time series original data

---

***Algorithm 1*** Determine appropriate model order $p, q, P, Q, D, d$ for SARIMA modeling approach.

---

For each high correlated feature list of (Branding type BT, Inventory size IS, Device category DC):

Construct Autocorrelation, and partial autocorrelation plots to heuristic determine an approximation of $p_0, q_0, P_0, Q_0, D_0, d_0$

Based on the previous step we calculate the AIC information criterion for $i\epsilon[p_0 - 2, p_0 + 2]$, $k\epsilon[q_0 - 2, q_0 + 2]$, $l\epsilon[P_0 - 2, P_0 + 2]$, $j\epsilon[Q_0 - 2, Q_0 + 2]$ and $D_0, d_0\epsilon\{0, 1\}$, where $l, j \leq p_0$, $q_0$ respectively.

Finally, we choose $p, q, P, Q, D, d$ that evaluate as $minAIC$ from the derived table.

---

# 4  Results

To test the software and the methodology that was created, all models have been applied on two different pricing rules over a specific website and the results of each analysis present revenue increase after providing the forecasted floor pricing values as new rules at Google Ad Manager pricing rules editing tool.
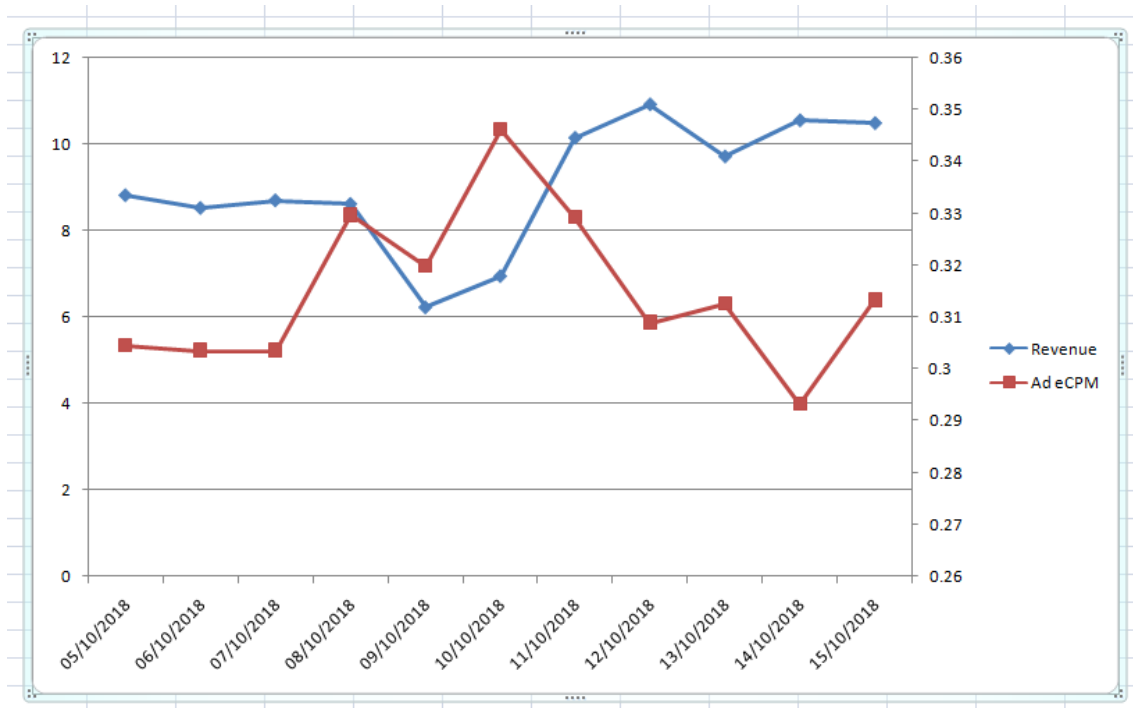


Figure 19: Forecasted values of Ad eCPM and their revenue

Figure 19 shows the revenue increment after applying the forecasted Ad eCPM values for 11 consequent days on the pricing rule floor price. Also, Table 5 presents the revenue uplift estimation for these days.

Table 5: Estimated revenue uplift case 1

| Dates | Estimated Revenue |
|---|---|
| 05-10-2018 | 8.82 € |
| 06-10-2018 | 8.53 € |
| 07-10-2018 | 8.70 € |
| 08-10-2018 | 8.62 € |
| 09-10-2018 | 6.23 € |
| 10-10-2018 | 6.96 € |
| 11-10-2018 | 10.14 € |
| 12-10-2018 | 10.91 € |
| 13-10-2018 | 9.71 € |
| 14-10-2018 | 10.55 € |
| 15-10-2018 | 10.48 € |

This revenue uplift was achieved at first by keeping the number of impressions nearly constant with load balancing of the available traffic of the specific ad units that are being examined. Only 20% of the real traffic appeared at the examined ad units were under the forecasting process.

Moreover, we acquired the resulting forecasted Ad eCPM values by the SARIMA process as depicted in Figure 17. Those values were applied as floor values for the specific pricing rule, and the resulting revenue increment was reported back as feedback from the Google Ad Manager reporting tool. This procedure of course must be at a constant level and should have daily feedback of the forecasting error so as to eliminate the case like Figure 20.

In Figure 20 on the other hand, shows that in this case although the forecasted Ad eCPM values presented uplift on the revenue for six days, the model error rate increased rapidly, and a re-evaluation should be instantiated. This fact led to a decrease in the estimated revenue which after a re-training and evaluation of the new forecasting values showed an increasing trend. Again, the uplift of the estimated revenue was increasing roughly higher than 10% for these six days as shown in Table 6.

Table 6: Estimated revenue uplift case 2

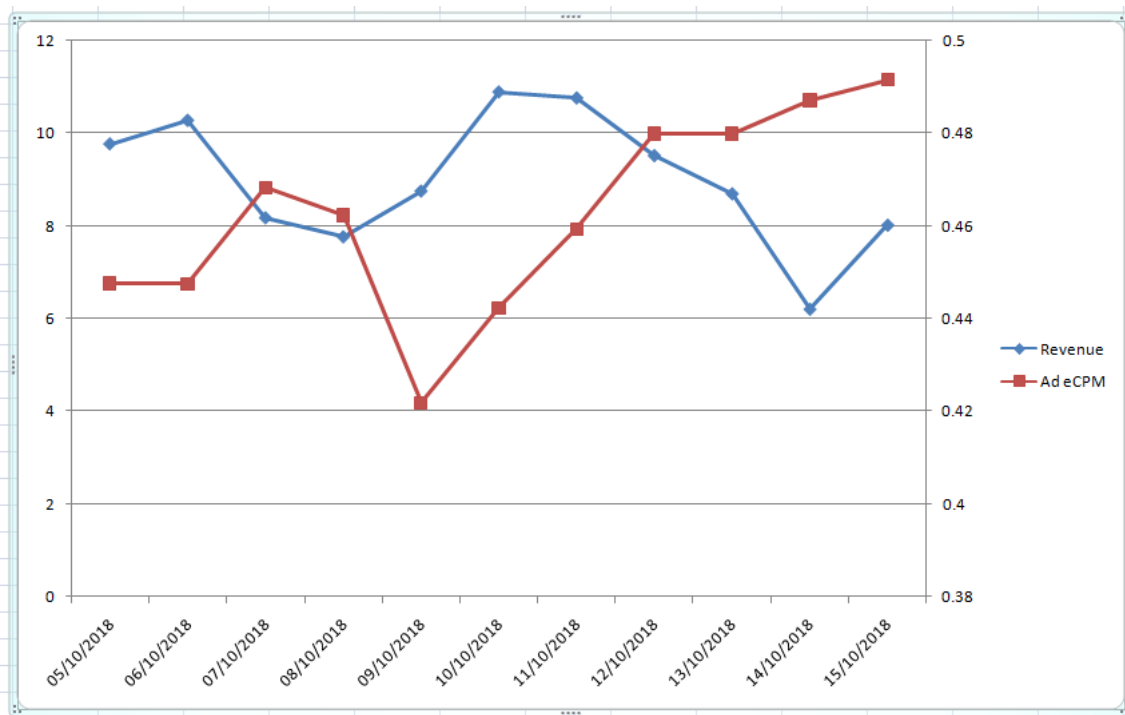| Dates | Estimated Revenue |
|---|---|
| 05-10-2018 | 9.78 € |
| 06-10-2018 | 10.28 € |
| 07-10-2018 | 8.18 € |
| 08-10-2018 | 7.77 € |
| 09-10-2018 | 8.76 € |
| 10-10-2018 | 10.90 € |



Figure 20: Forecasted values of Ad eCPM and their revenue

# 5  Software Usage

The software that has been developed is made on the Python 3.7.0 programming language. It has been used Jupyter notebook server for developing the test modeling approach locally as shown in Figure 22, while the final tool was deployed on a Google Colaboratory notebook.

Jupyter notebook environment installed as a local server can be obtained, is supported and fully documented on (Jupyter notebook environment , 2018)  as shown in Figure 23.

Google Collaboratory is a free Jupyter notebook environment that requires no setup and runs entirely in the cloud as displayed in Figure 21.

The advantage of Collaboratory over Jupyter notebook local server is that any machine learning and neural network algorithms can be accelerated using Google's Cloud computing engine and by using powerful graphic card GPUs or TensorFlow Processing Units - TPUs while the local server is limited to the computational resources provided by the personal computer.
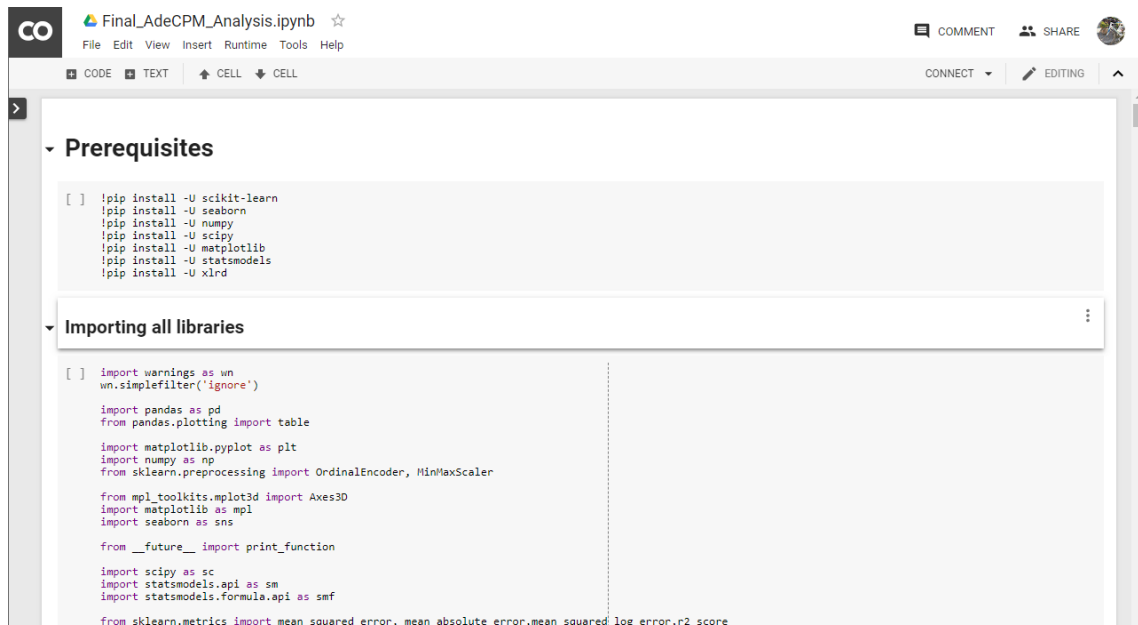


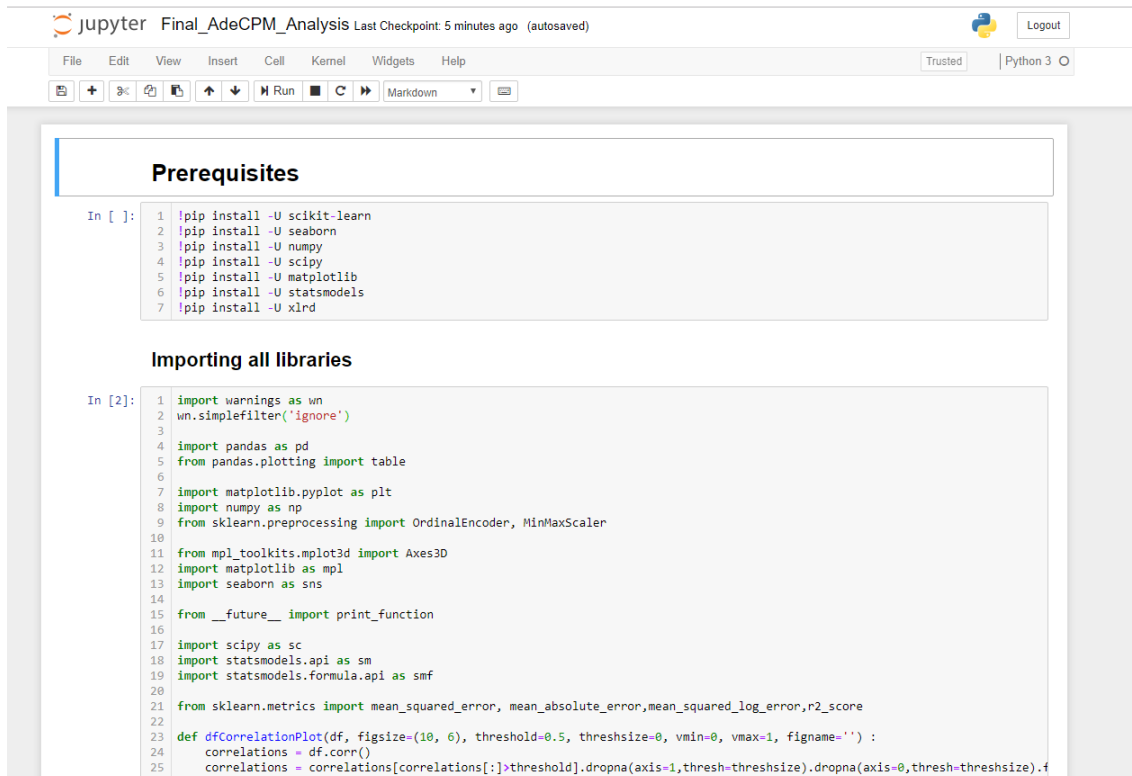Figure 21: Google Collaboratory online environment

Figure 22: Jupyter local notebook environment

In order to address the final implementation of the proposed software solution, certain popular libraries and toolboxes of Python were used such as (Numpy, 2018), (SciPy, 2018), (Pandas, 2018) and (Scikit-Learn, 2018) combined with two basic visualization libraries (MatplotLib, 2018) and (Seaborn, 2018).

For the statistical models, as well as for conducting statistical tests, and statistical data exploration such as ARIMA and SARIMA processes Statsmodels (Seabold, 2010) module has been used.

To visualize any interactivity, we also used (Bokeh, 2018) visualization library which supports such actions over plotted data.

Figure 23: Jupyter environment website

# 6  Conclusion

The problem this dissertation tried to solve was to create a time series algorithmic procedure such as to forecast floor pricing values of Google Ad Manager pricing rules, to achieve revenue increment (maximization) based on historical data. At first, the data was provided through a programmatic procedure based on the API of Google Ad Manager reporting tool and offered a dataset with several pricing rules over several websites with the same or different ad unit dimensions. The algorithm managed to work in multiple instances and is now becoming a pillar of how the Yield Management sector for publishers. Utilizing the important variables that affect pricing and trying to leverage the maximum available amount of advertising money is something all publishers need to do. With this algorithm and way of operating the way for maximum efficiency is being paved.

## 6.1  Future Steps

To make this procedure more independent and accurate on the forecasting values that are provided, we propose a deployment on a neural network environment where the decision and the application of each pricing value would be instant by the Artificial Intelligent algorithm. These kinds of algorithms eliminate the problem of real-time feedback and provide more efficient ways to evaluate forecasting models. Moreover, the ability to introduce external parameters other than the actual ad parameters, e.g., weather conditions, sudden viral news, acts of God is also something that we should heavily take into account. Since we have moved towards the first step of actually incorporating data to actively affect our decision to set a price, anything that is relevant to the ad will now become an important factor.

## 6.2  Weaknesses

The proposed approach also suffers of some weak points that focus on the feedback scheme that should instantiate a re-evaluation process phase as soon as the error rate of the forecasting values goes over a threshold value. Moreover, the time margin being only up to a day is not sufficient to go into a marginal analysis

# 7 References

*Bokeh*. (2018, December 1). Retrieved from Bokeh: https://bokeh.pydata.org/en/latest/

Bozdogan, H. (1987). Model Selection and Akaike's Information Criterion (AIC): The General Theory and Its Analytical Extensions. *Psychometrika* , pp. 345–370.

Jun Wang, S. Y. (2016). Real-Time Bidding based Display Advertising: Mechanisms and Algorithms. *ECIR* .

*Jupyter notebook environment* . (2018, December 1). Retrieved from Jupyter notebook environment : https://jupyter.org/

K.P. Burnham, D. A. (2004). Multimodel Inference: Understanding AIC and BIC in Model Selection. *Sociological methods and research* , pp. 261–304.

Kan Ren, W. Z. (2018, April). Bidding Machine: Learning to Bid for Directly Optimizing Profits in Display Advertising. *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING* .

*MatplotLib*. (2018, December 1). Retrieved from MatplotLib: https://matplotlib.org/

*Numpy*. (2018, December 1). Retrieved from Numpy 1.15.4: http://www.numpy.org/

*Pandas*. (2018, December 1). Retrieved from Pandas: http://pandas.pydata.org/

*Scikit-Learn*. (2018, December 1). Retrieved from Scikit-Learn: https://scikit-learn.org

*SciPy*. (2018, December 1). Retrieved from SciPy: https://scipy.org/

Seabold, S. a. (2010). Statsmodels: Econometric and statistical modeling with python. *9th Python in Science Conference* .

*Seaborn*. (2018, December 1). Retrieved from Seaborn: https://seaborn.pydata.org

Shuai Yuan, J. W. (2014). A Survey on Real Time Bidding Advertising. *IEEE ADKDD* .

Weinan Zhang, J. X. (2016). Learning, Prediction and Optimisation in RTB Display Advertising. *CIKM* .

Wush Chi-Hsuan Wu, M.-Y. Y.-S. (2015, August). Predicting Winning Price in Real Time Bidding with Censored Data. *ACM* .

# 8 Appendices

## 8.1 Software code

```
# -*- coding: utf-8 -*-
"""Final_AdeCPM_Analysis.ipynb

Automatically generated by Colaboratory.

Original file is located at
    https://colab.research.google.com/drive/1IuW6sKwpsGmyCSiDxD-fEoKAflENOOkL

#Prerequisites
"""

!pip install -U scikit-learn
!pip install -U seaborn
!pip install -U numpy
!pip install -U scipy
!pip install -U matplotlib
!pip install -U statsmodels
!pip install -U xlrd

"""##Importing all libraries"""

import warnings as wn
wn.simplefilter('ignore')

import pandas as pd
frompandas.plotting import table

importmatplotlib.pyplot as plt
importnumpy as np
```

[43]

```
33   fromsklearn.preprocessing import OrdinalEncoder, MinMaxScaler
34
35   from mpl_toolkits.mplot3d import Axes3D
36   importmatplotlib as mpl
37   importseaborn as sns
38
39   from __future__ import print_function
40
41   importscipy as sc
42   importstatsmodels.api as sm
43   importstatsmodels.formula.api as smf
44
45   fromsklearn.metrics                  import                  mean_squared_error,
46   mean_absolute_error,mean_squared_log_error,r2_score
47
48   defdfCorrelationPlot(df, figsize=(10, 6), threshold=0.5, threshsize=0, vmin=0, vmax=1, figname='') :
49   correlations = df.corr()
50   correlations                                                                    =
51   correlations[correlations[:]>threshold].dropna(axis=1,thresh=threshsize).dropna(axis=0,thresh=threshsiz
52   e).fillna(0)
53   fig = plt.figure(figsize=figsize)
54   ax = fig.add_subplot(111)
55   hm   =   sns.heatmap(round(correlations,3),   annot=True,   ax=ax,   cmap="coolwarm",fmt='.3f',
56   linewidths=.005, vmin=vmin, vmax=vmax)
57   if ~(figname=='') :
58   plt.savefig(figname)
59   plt.show()
60   return correlations
61
62   """# Data Import
63
64   ## Original Data
65
66   Importing Data from Excel spreadsheet.
67   """
68
69   originalDF = pd.read_excel('adwords_new_clean.xlsx', sheet_name='Data', index_col=None)
```

```python
70   prDF = pd.read_excel('pronews_price_floors.xlsx', sheet_name='DATA', index_col=[0,1,2])
71   originalDF = originalDF.join(prDF, on=['Pricing rules', 'Days', 'Branding types'])
72   print(originalDF.columns)
73   originalDF.describe()
74
75   """Some sample data rows."""
76
77   originalDF.head()
78
79   """## Encoding - Label indexing"""
80
81   feature_names = ['Pricing rules', 'Inventory sizes', 'Device categories', 'Branding types']
82   featuresDF = originalDF[feature_names]
83   OrdinalENC = OrdinalEncoder()
84   OrdinalENC.fit(featuresDF)
85   i=0
86   for cat in OrdinalENC.categories_:
87   print(feature_names[i], ' -> ', cat)
88       i+=1
89   featuresDFindexed   =   pd.DataFrame(OrdinalENC.transform(featuresDF),   columns=feature_names,
90   index=None)
91   df = originalDF[['Days', 'Ad requests', 'Ad impressions',
92               'Ad request eCPM (β,¬)', 'Ad eCPM (β,¬)', 'Floor',
93               'Diff AdCPMAdReqCPM', 'Estimated revenue (β,¬)']]
94   df[feature_names] = featuresDFindexed
95   df = df.set_index('Days')
96
97   """And the resulted dataframe."""
98
99   df.head()
100
101  """## Scaling
102
103  Scaling features that are too large to manipulate aside others.
104  """
105
106  MinMaxscaler = MinMaxScaler()
```

```
107    df[['Ad requests','Ad impressions']] = MinMaxscaler.fit_transform(df[['Ad requests','Ad impressions']])
108    df.head()
109
110    """## Basic Data-Set Plotting
111
112    ### 1. Floor value distribution per Pricing Rule
113    """
114
115    fromscipy.stats import norm
116    fig = plt.figure(figsize=(12, 6))
117    for i in range(1, len(OrdinalENC.categories_[0])+1):
118    ax = fig.add_subplot(2, 4, i)
119        ax.set_title(OrdinalENC.categories_[0][i-1].replace('pronews.gr', ''))
120    sns.distplot(df[df['Pricing rules']==i-1]['Floor'], ax=ax, fit=norm)
121    plt.tight_layout()
122    plt.show()
123
124    """### 2. Numerical features distribution"""
125
126    plt.figure(figsize=(12, 12))
127    g = sns.pairplot(df, hue='Pricing rules', vars=['Ad requests', 'Ad impressions', 'Ad eCPM (β,¬)', 'Floor',
128    'Estimated revenue (β,¬)'])
129    for i, j in zip(*np.triu_indices_from(g.axes, 1)):
130    g.axes[i, j].set_visible(False)
131
132    replacements = OrdinalENC.categories_[0]
133    for i in range(len(g.fig.get_children()[-1].texts)):
134    label = int(float(g.fig.get_children()[-1].texts[i].get_text()))
135        g.fig.get_children()[-1].texts[i].set_text(replacements[label].replace('pronews.gr ', ''))
136    g.fig.get_children()[-1].set_bbox_to_anchor((1.1, 0.5, 0, 0))
137    plt.show()
138
139    """### 3. Categorical features distribution"""
140
141    plt.figure(figsize=(12, 12))
142    g = sns.pairplot(df, hue='Pricing rules', vars=['Inventory sizes', 'Device categories', 'Branding types', 'Ad
143    eCPM (β,¬)'])
```

```python
144    for i, j in zip(*np.triu_indices_from(g.axes, 1)):
145    g.axes[i, j].set_visible(False)
146
147    replacements = OrdinalENC.categories_[0]
148    for i in range(len(g.fig.get_children()[-1].texts)):
149    label = int(float(g.fig.get_children()[-1].texts[i].get_text()))
150        g.fig.get_children()[-1].texts[i].set_text(replacements[label].replace('pronews.gr ', ''))
151    g.fig.get_children()[-1].set_bbox_to_anchor((1.1, 0.5, 0, 0))
152    plt.show()
153
154    """## Correlations"""
155
156    dfCorrelationPlot(df,threshold=0.0,threshsize=0,figname='CorrelationPlotDF.pdf')
157
158    """## Group by Inventory sizes vs Ad eCPM
159
160    Histogram plots that show the distribution of data against inventory sizes according to the Ad eCPM
161    values.
162    """
163
164    grpDaysInventory = df.reset_index().groupby(['Days','Inventory sizes'], as_index=False)['Ad eCPM
165    (β,¬)']
166    grpDaysInventory.aggregate(np.average)
167    DaysInventoryDF = grpDaysInventory.aggregate(np.average)[['Days', 'Inventory sizes','AdeCPM
168    (β,¬)']]
169    DaysInventoryDFpivot = DaysInventoryDF.pivot(index='Days', columns='Inventory sizes', values='Ad
170    eCPM (β,¬)').fillna(0)
171    axes = DaysInventoryDFpivot.hist(figsize=(16,20), layout=(5,5), sharey=True, sharex=True)
172    inventENC = OrdinalENC.categories_[1]
173    for i in range(axes.shape[0]):
174    for ax in axes[i]:
175    try:
176    num = int(float(ax.title.get_text()))
177    iflen(inventENC[num])>24:
178    inv = inventENC[num][:24]+'...'
179    else:
180    inv = inventENC[num]
```

[47]

```python
181    ax.set_title(inv+'\n (avg: '+str(round(DaysInventoryDFpivot[[num]].mean()[num], 2))+')')
182    ax.set_xlabel('Ad eCPM (β,¬)')
183    ax.set_xlim(0,0.4)
184    ax.set_ylabel('Frequency')
185    except:
186    continue
187    plt.show()
188
189    """## Group by Device categories vs Ad eCPM"""
190
191    grpDaysDevice = df.reset_index().groupby(['Days','Device categories'], as_index=False)['Ad eCPM
192    (β,¬)']
193    DaysDeviceDF = grpDaysDevice.aggregate(np.average)[['Days', 'Device categories','AdeCPM (β,¬)']]
194    DaysDeviceDFpivot = DaysDeviceDF.pivot(index='Days', columns='Device categories', values='Ad
195    eCPM (β,¬)').fillna(0)
196    axes = DaysDeviceDFpivot.hist(figsize=(9,8), layout=(2,2), sharey=True, sharex=True)
197    devicesENC = OrdinalENC.categories_[2]
198    for i in range(axes.shape[0]):
199    for ax in axes[i]:
200    num = int(float(ax.title.get_text()))
201    ax.set_title(devicesENC[num]+' (avg: '+str(round(DaysDeviceDFpivot[[num]].mean()[num], 2))+')')
202    ax.set_xlabel('Ad eCPM (β,¬)')
203    ax.set_xlim(0.15,0.4)
204    ax.set_ylabel('Frequency')
205    plt.show()
206
207    """## Group by Pricing rules vs Ad eCPM"""
208
209    grpDaysRules = df.reset_index().groupby(['Days','Pricing rules'], as_index=False)['Ad eCPM (β,¬)']
210    DaysRulesDF = grpDaysRules.aggregate(np.average)[['Days', 'Pricing rules','AdeCPM (β,¬)']]
211    DaysRulesDFpivot = DaysRulesDF.pivot(index='Days', columns='Pricing rules', values='Ad eCPM
212    (β,¬)').fillna(0)
213    axes = DaysRulesDFpivot.hist(figsize=(12,9), layout=(3,3), sharey=True, sharex=True)
214    devicesENC = OrdinalENC.categories_[0]
215    for i in range(axes.shape[0]):
216    for ax in axes[i]:
217    try:
```

[48]

```
218    num = int(float(ax.title.get_text()))
219    ax.set_title(devicesENC[num].replace('pronews.gr ', '')+
220                      ' (avg: '+str(round(DaysRulesDFpivot[[num]].mean()[num], 2))+')')
221    ax.set_xlabel('Ad eCPM (β,¬)')
222    ax.set_xlim(0,0.6)
223    ax.set_ylabel('Frequency')
224    except:
225    continue
226    plt.show()
227
228    """## Group by Branding types vs Ad eCPM"""
229
230    grpDaysBranding  =  df.reset_index().groupby(['Days','Branding  types'],  as_index=False)['Ad  eCPM
231    (β,¬)']
232    DaysBrandingDF = grpDaysBranding.aggregate(np.average)[['Days', 'Branding types','AdeCPM (β,¬)']]
233    DaysBrandingDFpivot = DaysBrandingDF.pivot(index='Days', columns='Branding types', values='Ad
234    eCPM (β,¬)').fillna(0)
235    axes = DaysBrandingDFpivot.hist(figsize=(10,4), layout=(1,2), sharey=True, sharex=True)
236    brandENC = OrdinalENC.categories_[3]
237    for i in range(axes.shape[0]):
238    for ax in axes[i]:
239    try:
240    num = int(float(ax.title.get_text()))
241    ax.set_title(brandENC[num]+' (avg: '+str(round(DaysBrandingDFpivot[[num]].mean()[num], 2))+')')
242    ax.set_xlabel('Ad eCPM (β,¬)')
243    ax.set_ylabel('Frequency')
244    except:
245    continue
246    plt.show()
247
248    """# Correlations Ad eCPM over Branding, Devices and Inventories
249
250    ##Correlation Matrix
251
252    Threshold value 0.95 and binding length at least 7
253    """
254
```

```python
255  newDF = pd.concat([DaysBrandingDFpivot, DaysDeviceDFpivot, DaysInventoryDFpivot],
256  axis=1, keys=['Branding','Device','Inventory'])
257  resCorrelations = dfCorrelationPlot(newDF,
258  threshold=0.95,
259  threshsize=7,
260  figsize=(12,8), vmin=0.92,
261  figname='CorrelationPlotBrandingDeviceInventory_0.95_7.pdf')
262
263  """##Correlated categorical dimensions
264  High scored features identification
265  """
266
267  resCorr = resCorrelations.reset_index()
268  for i in range(OrdinalENC.categories_[2].shape[0]):
269      r = resCorr
270  try:
271          r = r[(r['level_1']==i) & (r['level_0']=='Device')]
272          Branding                                                        =
273  int(r.Branding[r.Branding==np.max(r.Branding.as_matrix())].dropna(axis=1).columns[0])
274  inventory = int(r.Inventory[r.Inventory==np.max(r.Inventory.as_matrix())].dropna(axis=1).columns[0])
275  print(OrdinalENC.categories_[2][i])
276  print('\t'+str(Branding)+':'+OrdinalENC.categories_[3][Branding],
277  str(inventory)+':'+OrdinalENC.categories_[1][inventory])
278  except:
279  continue
280  print('=======================================')
281  for i in range(OrdinalENC.categories_[1].shape[0]):
282      r = resCorr
283  try:
284          r = r[(r['level_1']==i) & (r['level_0']=='Inventory')]
285          Branding                                                        =
286  int(r.Branding[r.Branding==np.max(r.Branding.as_matrix())].dropna(axis=1).columns[0])
287          Device = int(r.Device[r.Device==np.max(r.Device.as_matrix())].dropna(axis=1).columns[0])
288  print(i, OrdinalENC.categories_[1][i])
289  print('\t'+str(Branding)+':'+OrdinalENC.categories_[3][Branding],
290  str(Device)+':'+OrdinalENC.categories_[2][Device])
291  except:
```

```
292   continue
293
294   """#Model fitting phase ARIMA proccess
295
296   ##Analytical AIC criterion matrix
297
298   Analytical estimate of the Akaike information criterion (AIC) along with visualization of Q-Q plots per
299   model parameters.
300   """
301
302   cond = (df['Branding types']==1) &(df['Device categories']==1) &(df['Inventory sizes']==7)
303   data = df.loc[cond,['Ad eCPM (β,¬)']]
304   grpData = pd.DataFrame(data.groupby(['Days'], as_index=True)['Ad eCPM (β,¬)'].agg(np.average))
305   all_days = pd.date_range(grpData.index.min(), grpData.index.max(), freq='D')
306   grpData = grpData.reindex(all_days)
307   grpData = grpData.fillna(grpData.rolling(4,min_periods=1).mean())
308   grpData['Lag_1']=grpData - grpData.shift(1)
309   grpData['Lag_2']=grpData['Ad eCPM (β,¬)'] - grpData['Ad eCPM (β,¬)'].shift(2)
310   grpData['Lag_3']=grpData['Ad eCPM (β,¬)'] - grpData['Ad eCPM (β,¬)'].shift(3)
311   grpData['UP_DOWN']=(grpData['Lag_1']>=0).astype(int)
312   grpData = grpData.fillna(0)
313
314   maxp=3
315   maxd=0
316   maxq=2
317
318   aic_full = pd.DataFrame(np.zeros((maxp+1,maxq+1), dtype=float))
319   models = pd.DataFrame(np.zeros((maxp+1,maxq+1), dtype=object))
320
321   fig = plt.figure(figsize=(12,8))
322   ax = fig.add_subplot(111)
323
324   for p in np.arange(0,maxp+1):
325   for q in np.arange(0,maxq+1):
326   models.iloc[p,q] = sm.tsa.ARIMA(grpData['Ad eCPM (β,¬)'], order=(p,0,q))
327   try:
328   models.iloc[p,q] = models.iloc[p,q].fit(disp=False)
```

```python
329    aic_full.iloc[p,q] = models.iloc[p,q].aic
330            fig    =    sm.qqplot(models.iloc[p,q].resid,    line='q',    ax=ax,    fit=True,    label='Model
331    ('+str(p)+','+str(q)+')')
332    except:
333    aic_full.iloc[p,q] = 0.0
334    print(aic_full)
335
336    colormap = plt.cm.Pastel1
337    colors = [colormap(i) for i in np.linspace(0, 1, (maxp+1)*(maxq+1)*2)]
338    fori,j in enumerate(ax.lines):
339    j.set_color(colors[i])
340    plt.legend(loc='best')
341    plt.show()
342
343    """##Visualization of the models described in the previous step"""
344
345    fromdatetime import datetime as dt
346    fig = plt.figure(figsize=(12,8))
347    ax = fig.add_subplot(111)
348    grpData.loc[dt.strptime('2018-08-01         10:00:00','%Y-%m-%d         %H:%M:%S'):,['Ad         eCPM
349    (β,¬)']].plot(ax=ax)
350
351    d = 0
352
353    for p in np.arange(0,maxp+1):
354    for q in np.arange(0,maxq+1):
355    try:
356            predictions    =    models.iloc[p,q].predict(start=dt.strptime('2018-08-01    10:00:00','%Y-%m-%d
357    %H:%M:%S'),
358    end=dt.strptime('2018-10-10 10:00:00','%Y-%m-%d %H:%M:%S'),
359    dynamic=False)
360    predictions.shift(-1).plot(ax=ax, label='ARIMA('+str(p)+','+str(d)+','+str(q)+')')
361    except:
362    continue
363
364    colormap = plt.cm.Accent
365    colors = [colormap(i) for i in np.linspace(0, 1, (maxp+1)*(maxq+1)*2)]
```

```
366    fori,j in enumerate(ax.lines):
367    ifj.get_label()=='Ad eCPM (β,¬)':
368    j.set_color('r')
369    else:
370    j.set_color(colors[i])
371    plt.legend(loc='best')
372    plt.show()
373
374    """#Best model fit
375
376    Visualization of the predictions made by the model scored highest in the previous phase.
377    """
378
379    fromdatetime import datetime as dt
380    fig = plt.figure(figsize=(12,8))
381    ax = fig.add_subplot(111)
382    grpData.loc[dt.strptime('2018-07-10        10:00:00','%Y-%m-%d      %H:%M:%S'):,['Ad       eCPM
383    (β,¬)']].plot(ax=ax, color='r', marker='o')
384
385    p = 2
386    d = 0
387    q = 0
388    model = models.iloc[p,q]
389    predictions = model.predict(start=dt.strptime('2018-07-10 10:00:00','%Y-%m-%d %H:%M:%S'),
390    end=dt.strptime('2018-10-31 10:00:00','%Y-%m-%d %H:%M:%S'),
391    dynamic=False)
392
393    predictions.shift(-1).plot(ax=ax, label='ARIMA('+str(p)+','+str(d)+','+str(q)+')', color='b', marker='o')
394
395    test = grpData.loc[dt.strptime('2018-07-10 10:00:00','%Y-%m-%d %H:%M:%S'):,['Ad eCPM (β,¬)']]
396    fromsklearn.metrics                          import                    mean_squared_error,
397    mean_absolute_error,mean_squared_log_error,r2_score
398
399    tbl = pd.DataFrame({
400    'AIC' :models.iloc[p,q].aic,
401    'BIC' :models.iloc[p,q].bic,
402    'r2' : r2_score(test, predictions.iloc[:92]),
```

[53]

```
403    'MAE' :mean_absolute_error(test, predictions.iloc[:92]),
404    'MSE' :mean_squared_error(test, predictions.iloc[:92]),
405    'RMSE' :np.sqrt(mean_squared_error(test, predictions.iloc[:92]))
406    }, index=['ARIMA('+str(p)+','+str(d)+','+str(q)+')'])
407
408    table(ax, np.round(tbl.T, 4), loc='center right', colWidths=[0.2, 0.2])
409
410    plt.legend(loc='best')
411    plt.show()
412
413    """##Model Summary"""
414
415    print(model.summary2())
416
417    """##Forecasted values
418
419    The forecasted values of the previous model.
420    """
421
422    print('Forecasts')
423    print(predictions.iloc[92:])
424
425    """#Seasonal ARIMA proccess
426
427    ##Data decomposition
428
429    Decomposing data to show thier trend, seasonality
430    """
431
432    fromstatsmodels.tsa.seasonal import seasonal_decompose
433    cond = (df['Branding types']==1) &(df['Device categories']==1) &(df['Inventory sizes']==7)
434    data = df.loc[cond,['Ad eCPM (β,¬)']]
435    grpData = pd.DataFrame(data.groupby(['Days'], as_index=True)['Ad eCPM (β,¬)'].agg(np.average))
436    all_days = pd.date_range(grpData.index.min(), grpData.index.max(), freq='D')
437    grpData = grpData.reindex(all_days)
438    grpData = grpData.fillna(grpData.rolling(4,min_periods=1).mean())
439
```

```python
440    ts = grpData.loc[dt.strptime('2018-07-01 10:00:00','%Y-%m-%d %H:%M:%S'):,['Ad eCPM (β,¬)']]
441    decomp = seasonal_decompose(ts)
442
443    tr = decomp.trend
444    ses = decomp.seasonal
445
446    plt.figure(figsize=(12,8))
447    plt.subplot(311)
448    plt.plot(ts, label='Original', marker='o')
449    plt.legend(loc='best')
450    plt.subplot(312)
451    plt.plot(tr, label='trend', marker='o', color='r')
452    plt.legend(loc='best')
453    plt.subplot(313)
454    plt.plot(ses, label='seasonal', marker='o', color='g')
455    plt.legend(loc='best')
456    plt.show()
457
458    """##ACF and PACF plots
459
460    Also show their autocorrelation and partial autocorrelation plots.
461    """
462
463    fig = plt.figure(figsize=(12,7))
464    ax1 = fig.add_subplot(211)
465    fig = sm.graphics.tsa.plot_acf(ts.values.squeeze(), lags=40, ax=ax1)
466    ax2 = fig.add_subplot(212)
467    fig = sm.graphics.tsa.plot_pacf(ts, lags=40, ax=ax2)
468
469    """## SARIMA model fitting"""
470
471    p = 3
472    d = 0
473    q = 1
474    sesP = 1
475    sesD = 0
476    sesQ = 0
```

```
477    ses = 2
478
479    nm       =       sm.tsa.statespace.SARIMAX(grpData['Ad       eCPM       (β,¬)'],       order=(p,d,q),
480    seasonal_order=(sesP,sesD,sesQ,ses)).fit( disp=False)
481    predictions = nm.predict(start='2018-07-10 10:00:00', end='2018-10-09 10:00:00', dynamic=False)
482    print(nm.summary())
483
484    """##Visaulizing data"""
485
486    fig, ax = plt.subplots(figsize=(12,8))
487    grpData.loc[dt.strptime('2018-07-10       10:00:00','%Y-%m-%d       %H:%M:%S'):,['Ad       eCPM
488    (β,¬)']].plot(ax=ax,color='b',
489    label='Original Observations')
490    predictions.shift(-1).plot(ax=ax,
491    label='SARIMAX('+str(p)+','+str(d)+','+str(q)+')x('+str(sesP)+','+str(sesD)+','+str(sesQ)+')
492    m='+str(ses),color='r')
493
494    test = grpData.loc[dt.strptime('2018-07-10 10:00:00','%Y-%m-%d %H:%M:%S'):,['Ad eCPM (β,¬)']]
495
496    tbl = pd.DataFrame({
497    'AIC' :nm.aic,
498    'BIC' :nm.bic,
499    'r2' : r2_score(test, predictions),
500    'MAE' :mean_absolute_error(test, predictions),
501    'MSE' :mean_squared_error(test, predictions),
502    'RMSE' :np.sqrt(mean_squared_error(test, predictions))
503    }, index=['Model('+str(p)+','+str(q)+')'])
504
505    table(ax, np.round(tbl.T, 4), loc='center right', colWidths=[0.2, 0.2])
506
507    plt.legend(loc='best')
508    plt.show()
509
510    """##Forecasting 7 steps ahead"""
511
512    print(nm.forecast(steps=7))
513
```

```python
514    """##Plotting residuals"""
515
516    residuals = pd.DataFrame(nm.resid, columns=['Residuals'])
517    ax = residuals.plot.hist(figsize=(10,6))
518    residuals.Residuals.plot.kde(ax=ax, secondary_y=True)
519    importscipy.stats as sc
520    x=np.arange(-0.15,0.15,0.001)
521    plt.plot(x,sc.norm.pdf(x,0,0.0265),label='N(0,1)')
522    plt.legend(loc='upper left')
523    plt.show()
524
525    """#Last method on modeling data
526
527    ##Exponential Smoothing
528
529    Specifically Holt - Winters method on Exponential Smoothing
530    """
531
532    fromstatsmodels.tsa.holtwinters import ExponentialSmoothing
533    Holtdf = grpData['Ad eCPM (β,¬)']
534    split = round(len(grpData['Ad eCPM (β,¬)'])*0.8)
535    train, test = Holtdf.iloc[:split], Holtdf.iloc[split:]
536    model = ExponentialSmoothing(train, seasonal_periods=9, seasonal="mul").fit()
537    pred = model.predict(start=test.index[0], end=test.index[-1])
538
539    fig, ax = plt.subplots(figsize=(10,6))
540    plt.plot(train.index, train, label='Train')
541    plt.plot(test.index, test, label='Test')
542    plt.plot(pred.index, pred, label='Holt-Winters')
543
544    tbl = pd.DataFrame({
545    'AIC' :model.aic,
546    'BIC' :model.bic,
547    'r2' : r2_score(test, pred),
548    'MAE' :mean_absolute_error(test, pred),
549    'MSE' :mean_squared_error(test, pred),
550    'RMSE' :np.sqrt(mean_squared_error(test, pred))
```

```
551    }, index=['Holt-Winters'])
552
553    table(ax, np.round(tbl.T, 4), loc='center right', colWidths=[0.2, 0.2])
554    plt.legend(loc='best')
555    plt.show()
556
557    """##Residuals of Holt - Winters method"""
558
559    residuals = pd.DataFrame(model.resid, columns=['Residuals'])
560    ax = residuals.plot.hist()
561    residuals.Residuals.plot.kde(ax=ax, secondary_y=True, color='b',label='kde')
562
563    importscipy.stats as sc
564    x=np.arange(-0.15,0.15,0.001)
565    plt.plot(x,sc.norm.pdf(x,0,0.02),label='N(0,1)')
566
567    plt.legend(loc='upper left')
568    plt.show()
```