

Google Trends and Tourists' Arrivals: Emerging biases and proposed corrections

Theologos Dergiades[‡]

Department of International & European Studies
University of Macedonia, Greece
e-mail: dergiades@uom.edu.gr

Eleni Mavragani

School of Economics, Business Administration and
Legal Studies
International Hellenic University
e-mail: e.mavragani@ihu.edu.gr

and

Bing Pan

Department of Recreation, Park & Tourism Management
Pennsylvania State University
e-mail: bingpan@psu.edu

Abstract

As search engines constitute a leading tool in planning vacations, researchers have adopted search engine query data to predict the consumption of tourism products. However, when the prevailing shares of visitors come from countries in different languages and with different dominating search engine platforms, the identification of the aggregate search intensity index to forecast overall international arrivals, becomes challenging since two critical sources of bias are involved. After defining the language bias and the platform bias, this study focuses on a destination with a multilingual set of source markets along with different dominating search engine platforms. We analyze monthly data (2004-2015) for Cyprus with two non-causality testing procedures. We find that the corrected aggregate search engine volume index, adjusted for different search languages and different search platforms, is preferable in forecasting international visitor volumes compared to the non-adjusted index.

Keywords: Web search intensity; Google Trends; Tourists' arrivals.

[‡] Corresponding author; Department of International & European Studies, University of Macedonia, PC 546 36, Egnantia Str., Thessaloniki Greece.

1. Introduction

Over the recent years, the availability of data gleaned from copious web sources (social media, search engines, etc.) sparked a new interest in the area named real-time economics. In one of the earliest studies, Choi and Varian (2009) demonstrated that properly selected query indices provided by Google are useful in forecasting activities in different economic sectors, such as the automobile industry and the tourism market. Their study has triggered a flurry of scientific publications that use web-related data, which aim to explain upcoming trends in various markets, including foreign exchange markets, stock markets, sovereign bond markets, labor markets or even real estate markets (see among others Joseph *et al.*, 2011; Smith, 2012; Beracha and Wintoki 2013; Dergiades *et al.*, 2015). Credible evidence shows that web-related data offer added value when it comes to predicting upcoming economic activities.

Forecasting tourism demand is essential for practitioners and policymakers. Accurate forecasts provide valuable aid for a) the development of medium- to long-run marketing and tourism strategies, b) the formation of pricing policies, c) the appropriate scheduling of investments (Clerides and Adamou, 2010), and d) the effective allocation of the limited resources (Song *et al.*, 2006; Yang *et al.*, 2015). Nowadays, web search engines constitute one major tool in planning vacations and can help to improve demand forecasting for the tourism product. In this study, we argue that the failure to account for two sources of bias (*language bias* and *platform bias*) frequently encountered during the construction of *Search Intensity Indices (SII)* from search engines, deteriorates the quality of the delivered index as a predictor.

We argue that a *SII* based on search engines in one language is unbiased, only if all the visitors perform their web searches in one language. In more details, as we use the search query volume in one language (e.g., English), the web search intensity is attributed just to a set of countries (U.S.A., U.K., etc.), while neglecting entirely the web search intensity formed in other languages. Thus, failure to account for all the languages that correspond to the respective source markets will give rise to the first source of bias, *language bias*. In addition, to protect the privacy of search engine users, the dominating search platform Google does not deliver data if the search volume for certain keywords is relatively small. Consequently, one cannot construct an entirely accurate aggregate index if some international tourists who searched on Google speak a rare language. One can imagine most countries will have a small number of international arrivals speak rare languages. Hence, this *language bias* is not a question of presence or absence, but rather an

existing problem in various degrees. Even if at some point in our sample, all primary source markets use the same language, there is no guarantee that this will be the case in the future.

A second bias may exist if the search engine used to collect data is not the only platform in the source market of interest - thus, the *platform bias*. In such cases, the measured volume of queries underestimates the actual volume of relevant queries (the search volume from other engines is ignored), failing to convey the precise interest of users and its evolution over time.

This study concentrates on Cyprus and evaluates the impact of the relevant web *SII*, captured by search platforms, on the consumption of the tourism product. Cyprus is an ideal candidate country since the composition of international arrivals makes both sources of bias coexist. It allows us to examine how we can deal with the effects of the *language bias* and the *platform bias*, with a purpose to construct an effective predictor for international arrivals. We concentrate on the search engine of Google for two main reasons: Google is the most popular search engine globally, with a market share amounting to 66.7% (Yang *et al.*, 2015); Google provides the historical intensity of the conducted queries through a platform called Google Trends (<https://www.google.gr/trends>).

Accurate prediction of the international arrivals in Cyprus is crucial since the overall contribution of the tourism industry in 2014 is more than €3 billion, a 21.3% of the GDP (KPMG, April 2016, available at: <https://www.kpmg.com/cy/>). Projections for the next ten years show that the absolute contribution of the tourism industry is expected to grow at a steady annual rate of around 5%. By 2025, the relative contribution of the tourism sector is anticipated to reach 25.5%. In addition, only around 40% of international arrivals are from English-speaking countries in 2015.¹ Around 30% of visitors speaks Russian, Greek, German, and Swedish as their native languages. Thus, English keyword searches might not represent a majority of searches for the country. Furthermore, Google is not the dominant search engine in the Russian market. A search engine called Yandex on average operates approximately 60% of the Russian market, while Google's respective share is about 25% (see www.liveinternet.ru).

This study adopts two non-causality testing techniques, in the time domain and the frequency domain. It introduces an uncomplicated way to select appropriate keywords, and investigates the predictive power of Google's *SII* towards the arrivals of international tourists in Cyprus at an aggregate and disaggregate level. The findings show that the

¹To the best of our knowledge the only study that deals with a destination that receives visitors from different countries is that of Choi and Varian (2012). Choi and Varian (2012) act at a disaggregated level only, and they do not provide many details about the construction of the search intensity index (e.g., keywords used).

presence of the *language bias*, and the *platform bias* render the simple aggregate *SII* ineffective in predicting the total number of international arrivals. The corrected aggregate *SII* conveys a more valuable predictive content.

Our study has the following structure: Section 2 briefly reviews the literature devoted to the broad field of econometric forecasting through web-related data, paying particular attention to the tourism market. Section 3 illustrates the methodological framework and section 4 presents the data and the preliminary econometric analysis. Section 5 presents our main empirical findings while the resulting managerial implications are discussed in Section 6. Finally, Section 7 concludes this study.

2. Literature Review

Researchers try to provide accurate forecasts for the arrivals of tourists implementing a wide range of techniques. Peng *et al.* (2014) summarize two broad categories of techniques: time-series econometrics and artificial intelligence methods. The former category includes econometric models ranging from very simple univariate specifications (Geurts and Ibrahim, 1975; Martin and Witt, 1989) to more advanced multivariate specifications (Halicioglu, 2010; or Bangwayo-Skeete and Skeete, 2015). The latter category comprises models ranging from artificial neural networks (Burger *et al.*, 2001) to genetic algorithms (see among others, Chen and Wang, 2007). A detailed review on the topic is discussed in Peng *et al.* (2014) and Song *et al.* (2003).

The empirical studies on tourism demand introduce an extensive set of explanatory factors to model arrivals. Using a diverse set of criteria, several researchers have grouped these factors (see Frechtling, 2001; Middleton *et al.*, 2009). Frechtling (2001) groups tourism demand factors into: 1) push, 2) pull, and 3) resistance factors. All groups above embrace both quantitative and qualitative factors, with the former to be those most frequently used in the empirical analysis since they are easily measurable and accessible effortless. In contrast, while qualitative attributes play a very crucial role in determining arrivals, rarely are these incorporated in demand specifications as their quantification is an arduous task.

In more detail, push factors include features related to the source markets. For example, Martins *et al.* (2017) find that the per capita income is critical in explaining arrivals (based on a large panel of 218 countries) while Goh *et al.* (2008) shows, for the case of Hong Kong, that leisure time (in the two sources markets - U.S.A. and U.K.), influences

arrivals stronger than economic factors. Additionally, Dragouni *et al.* (2016), focusing on the U.S. outbound tourism, support that the effect of consumers' sentiment and mood on the demand for tourism appears significant but time and event dependent.

Pull factors refer to attributes of the destination country (the quality the natural resources, Foreign Direct Investments (FDI) or social ties, etc.). For instance, Deng *et al.* (2002) mention that natural resources constitute one of the leading attractions for tourism demand. A report of the World Tourism Organization (UNWTO) in 2012, evaluates that the number of travelers attracted by natural resources is predicted to rise rapidly over the upcoming decades, at a rate higher than the average of the tourism industry. Furthermore, by surveying the attitude of 2,356 individuals from Italy with respect to the exploitation of natural resources, Meleddu and Pulina (2016) identify a positive propensity to pay a premium for eco-tourism. Moreover, Foreign Direct Investments may increase the overall quality of the provided services in the tourism sector from several aspects. Hence, increase in the number of arrivals is expected. The above is verified by Tang *et al.* (2007) for China and Craigwell and Winston (2008) for 21 small island countries. Gafer and Tchetchik (2017), focusing on Israel, recognize the significance of social ties in affecting tourism demand.

Finally, resistance factors refer to features that constrain traveling among the source markets and the destination. For instance, Turner and Witt (2001) find that for New Zealand, relative prices are significant in explaining arrivals. A recent study by Poprawe (2015), using a panel dataset for more than 100 countries, estimates that if the perceived corruption decreases by one unit, then tourists' arrivals increase in the range of 2% to 7%. The adverse effects of corruption in tourism demand are also confirmed by the studies of Saha and Yap (2015) and Das and Dirienzo (2010), using a panel of 130 and 119 countries, respectively.

Apart from the factors discussed above, the increasing availability of data capturing consumers' online activities has led many studies to adopt the identified web search activity as a tourism demand predictor. Yang *et al.*, (2014) recognize that the major advantages of such data lie in: a) reveal preferences in real-time, b) provide data in relatively high frequency (e.g., daily or weekly) and, c) depict changes in consumers' preferences. The latter advantage consists a solution to the inherent specification problem often encountered in traditional univariate time-series models (e.g., ARMA models). Traditional models fail to provide robust forecasts when sudden one-off events are taking place and alter the pattern of the series. Of course, empirical applications using web-related data have

been conducted for several markets besides the tourism market. For instance, Smith (2012) shows that the online search intensity, as captured by Google, significantly explains movements in the currency markets. Joseph *et al.*, (2011), using data from the Google Trends, forecast abnormal stock returns and the respective trading volume for the respective stock tickers. Based on a sample of 3000 stocks, Da *et al.*, (2011) argue that a higher search volume index for the relevant stock ticker forecasts higher stock prices in the short-run. Beracha and Wintoki (2013) find that the abnormal search intensity in the real estate market of a city predicts the abnormal housing prices. Finally, Dergiades *et al.*, (2015) show that the web search intensity for the keyword *Grexist*, explains future price movements of the 10-year government Greek bonds.

A substantial number of web users seek information through search engines before taking a trip (Fesenmaier *et al.*, 2011). Despite the large volume of studies dedicated to forecasting the demand of the tourism product, there is relatively a small number of studies adopting web search intensity data. Xiang and Pan (2011) analyze search queries of U.S. cities and find that “the ratio of travel queries among all queries about a specific city seems to associate with the touristic level of that city” (p. 88). Choi and Varian (2012) validate that search intensity data provided by Google for nine source markets-countries are indeed useful predictors of tourists’ arrivals to Hong Kong from each respective market.

Yang *et al.* (2015) implemented an ARMA -Autoregressive Moving Average-specification and the standard Granger non-causality test and affirmed that query volume data from two search engines - Google and Baidu - contribute significantly to decreasing forecasting errors when predicting the number of visitors to Hainan (a Chinese province). Bangwayo-Skeete and Skeete (2015) direct their interest to international visitors to five Caribbean destinations (Jamaica, Bahamas, Dominican Republic, Cayman, and St. Lucia). They conduct their analysis by implementing a simple AR-MIDAS model, a SARIMA model (Seasonal Autoregressive Integrated Moving Average), and a benchmark AR model (Autoregressive). The former model appears to perform the best in most of the conducted pseudo-forecasting experiments. Overall, the authors argue that after the proper construction of the Google search intensity indicator, significant gains are achieved in forecasting tourist arrivals.

These studies validated the value of search engine intensity data in predicting tourist arrivals. However, in most of these studies, the dominant visitor source markets are English-speaking countries; and almost all the source markets use Google as the dominant search engine. For instance, in the study by Bangwayo-Skeete and Skeete (2015), the three

source markets for the five investigated destinations are U.S.A., U.K. and Canada with the market share of Google's search engine in those countries to be 68.8%, 92.7% and 92.9%, respectively (Kennedy and Hauksson, 2012). However, the source market in many countries may use a variety of languages, and Google might not be the dominant search engine. For example, Canada, among several other countries, is officially a bilingual country or in numerous other countries, large linguistic minorities exist. Moreover, Google is not the dominant search engine in some large markets such as Russia and China, where Yandex and Baidu have a market share approximately equal to 60% and 52%, respectively (Kennedy and Hauksson, 2012). In addition, if we use only one language in the search engine or if we focus on one search engine with a small market share (when another search engine is the market leader), then we are neglecting the web search intensity formed in other languages and other search engines.

This study adopted Cyprus as a case study, where among its source markets, English is not the dominant language; and some source markets do not use Google as the major search engine. We are interested in finding out how to acquire search engine query data in forecasting international tourist arrivals in this country.

3. Methodology

The first approach implemented to examine the hypothesis of no predictability is the standard Granger causality test (Granger, 1969) using a VAR specification (Sims, 1980). The null hypothesis is examined by testing whether lagged values of one variable may significantly contribute in predicting current values of another variable. Additionally, as a second testing procedure, we employ the Breitung and Candelon (2006) (B&C, hereafter) non-causality test since it illustrates features that cannot be traced in the standard Granger causality test. For instance, a) it permits to distinguish the dynamic characteristics of the causal relationship, either short-run or long-run causality, b) it can identify causal relationships even if the underlying linkage among the variables of interest is non-linear, and finally, c) it delivers robust results in the presence of volatility clusters, a common characteristic of data with high frequency (Breitung and Candelon, 2006).

B&C propose a procedure to test for non-causality at the frequency domain by exploiting the Cholesky structural representation of a VAR model. Once a VAR model is estimated for the $\mathbf{Z}_t = (A_t \ G_t)^T$ 2×1 vector of stationary variables, from the structural

representation of the model, the predictive content of G_t towards A_t can be tested through the following Fourier transformation that takes place on the moving average coefficients:

$$M_{G \rightarrow A}(\omega) = \log \left[1 + \frac{|\Psi_{12}(e^{-i\omega})|^2}{|\Psi_{11}(e^{-i\omega})|^2} \right] \quad (1)$$

If G_t does not cause A_t at frequency ω , $|\Psi_{12}(e^{-i\omega})|^2$ should be equal to zero. Given that $|\Psi_{12}(e^{-i\omega})|^2$ is a complicated non-linear function, B&C propose to test the same hypothesis using the following set of restrictions.²

$$\sum_{k=1}^p \theta_{12,k} \cos(k\omega) = 0 \quad \text{and} \quad \sum_{k=1}^p \theta_{12,k} \sin(k\omega) = 0 \quad (2)$$

B&C examine the validity of the linear restrictions illustrated in *eq. (2)*, for frequencies ω that receive values within the interval of $(0, \pi)$, by comparing the estimated statistic with the 0.05 critical value of the χ^2 distribution with 2 degrees of freedom.

4. Data and Preliminary Econometric Analysis

4.1. Data Sources

This study employs monthly time-series data on tourist arrivals in Cyprus along with web search intensity data for appropriately selected keywords. The range of our sample is January 2004 to April 2016 (148 observations) due to data availability from Google Trends. The data for the total arrivals of tourists in Cyprus (see Fig. 1) as well as arrivals for each origin country come from the Statistical Service of Cyprus (Fig. 2 shows the arrivals per country as a market share). The arrivals per country are available until December 2015 (144 observations). For the selected sample, to extract the *SII* related to the tourist product of Cyprus, we use the Google Trends facility.

² Given that $\sin(k\omega) = 0$ in the cases where $\omega = 0$ and $\omega = \pi$, then it comes that the second restriction in *eq. (2)* is simply disregarded.

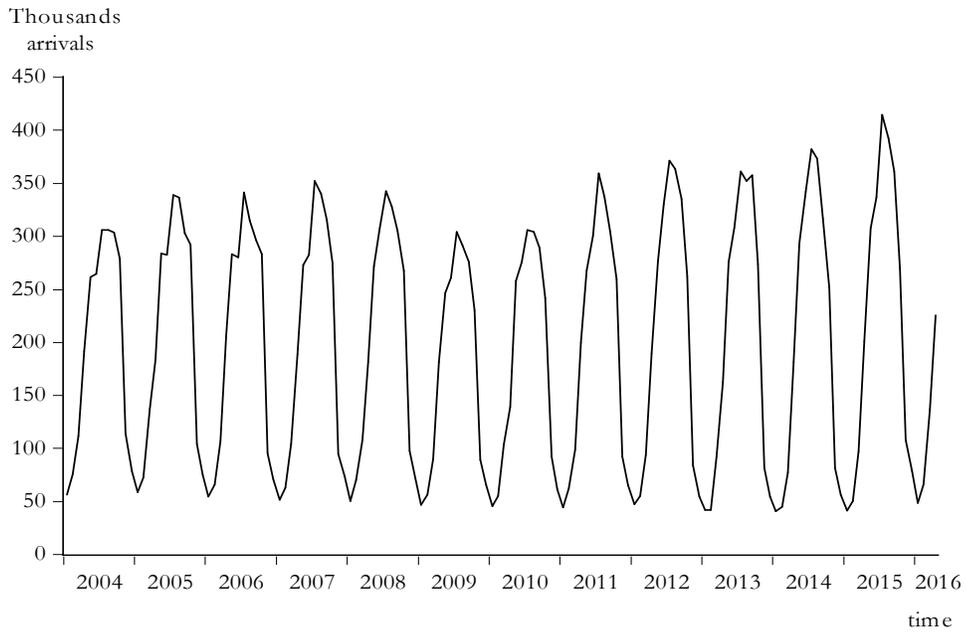


Figure 1. Monthly arrivals of tourists in Cyprus

To capture the entire web search intensity for a destination, we first consider possible existence of the *language bias* and the *platform bias*. To identify the presence of the first source of bias, we disentangle the aggregate number of tourist arrivals in Cyprus by country of origin with a purpose to specify the corresponding languages. The market share in the total arrivals per country is illustrated in Fig. 2. Visual inspection suggests that five countries are the main source markets, representing jointly 74.1% of the market share (average share of the total monthly arrivals during the period of study 2004-2015), while the respective share for all the other countries is 25.9%. The major source markets countries are the following: UK (45.4%), Russia (10.1%), Greece (7.7%), Germany (7.2%) and Sweden (3.7%). Hence, we concentrate on the respective languages of English, Russian, Greek, German and Swedish.

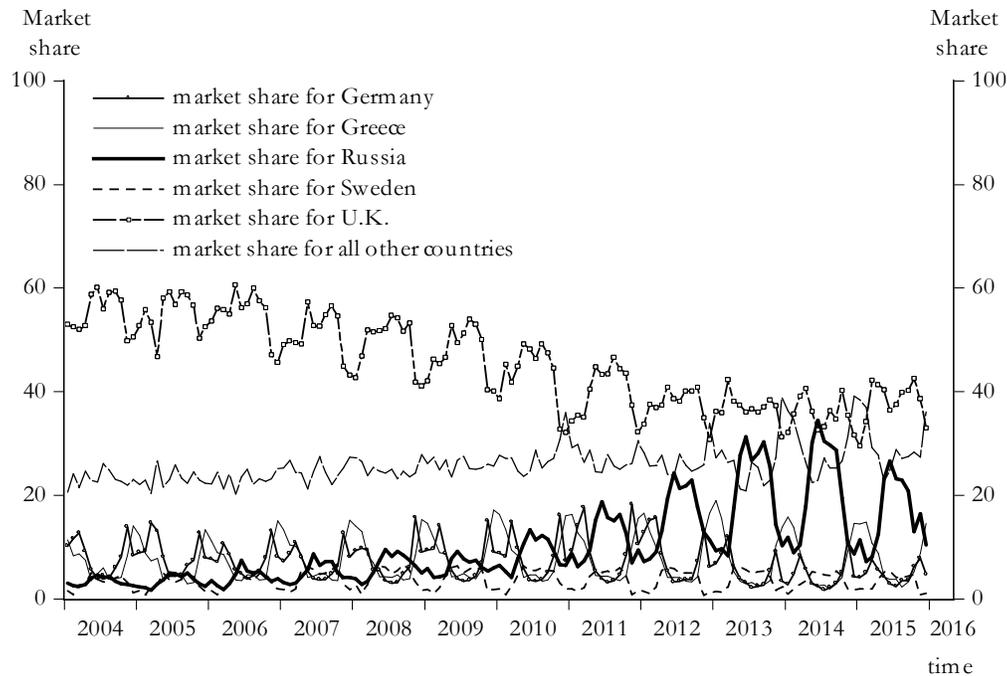


Figure 2. Markets shares in tourists' arrivals for Cyprus per country

Fig. 2 shows that the market shares for each country evolve differently. The market share of the U.K. shrinks from 55.2% in 2004 to 38.0% in 2015. The German share in 2004 was 7.2%, reaching the value of 4.9% in 2015. The Russian market share increased from 3.3% in 2004 to 16.3% at the end of the sample. The total market share of the rest of countries rose from 23.2% in 2004 to 29.6% in 2015. Finally, the shares for Greece and Sweden have stayed relatively constant.

The next step is to identify the appropriate keywords that are directly related to a potential visit to Cyprus for each language. It is reasonable to assume that the performed searches contain the term *Cyprus*. The identification of appropriate keywords involves the following steps: 1) by first selecting the source market of interest, we type the term *Cyprus* (in the language of the source market) in the Google Correlate tool to attain other queries that present similar patterns (the similarity is ascertained through a simple correlation coefficient). From the delivered queries, ranked in terms of correlation, we select the query that presents the highest correlation to our search term and its meaning refers explicitly to a visit in Cyprus (e.g., *flights to Cyprus*). 2) From Google Trends facility, we extract a monthly frequency series for the keyword identified in the previous step. We verify the validity of our chosen keyword by examining the top related queries as suggested by Google Trends. If the vast majority of the related queries imply interest for visiting Cyprus, we may argue

in favor of our keyword. 3) For those cases where in the first step our initial key term (e.g., *Cyprus*) does not deliver keywords that convey direct interest for a trip to the destination, we type our keyword (*Cyprus*) in Google Trends facility and from the delivered related queries, we select the one that expresses explicit interest to visit the destination.

For instance, in the case of U.K., the Google Correlate facility suggests that the most highly correlated term to *Cyprus* (which implies explicit intention to visit Cyprus) is the keyword *hotel Cyprus*. At the second stage, we type *hotel Cyprus* in the Google Trends facility, and we examine the relevant queries. All the relevant queries (hotel in Cyprus, Paphos Cyprus, Paphos, hotels Cyprus, Cyprus holidays, Portaras Cyprus, Portaras) verify the validity of our selected keyword since they imply direct interest to visit Cyprus. The finally extracted index in monthly frequency is presented in Fig. 3a. Implementing the same strategy for the remaining source markets, we end up with the following keywords. For the Russian market, the identified keyword is туры кипр (tours Cyprus), and the respective index is illustrated in Fig. 3b. For the German market, the keyword is *hotel zypern* (hotel Cyprus) and depicted in Fig. 3c, and for the Swedish market, the keyword is *cypern resor* (Cyprus travel) shown in Fig. 3d. Finally, the strategy failed to deliver a keyword that expresses an intention to visit Cyprus for the case of Greece. We tried keywords that are similar to those identified for the other countries, as for example *ξενοδοχεία Κύπρος* (hotels Cyprus) or *διακοπές Κύπρος* (holidays Cyprus), and the Google Trends facility indicated that there is not enough search volume to deliver results. Therefore, we are unable to construct a web *SII* for Greece, and we proceed with the remaining markets.

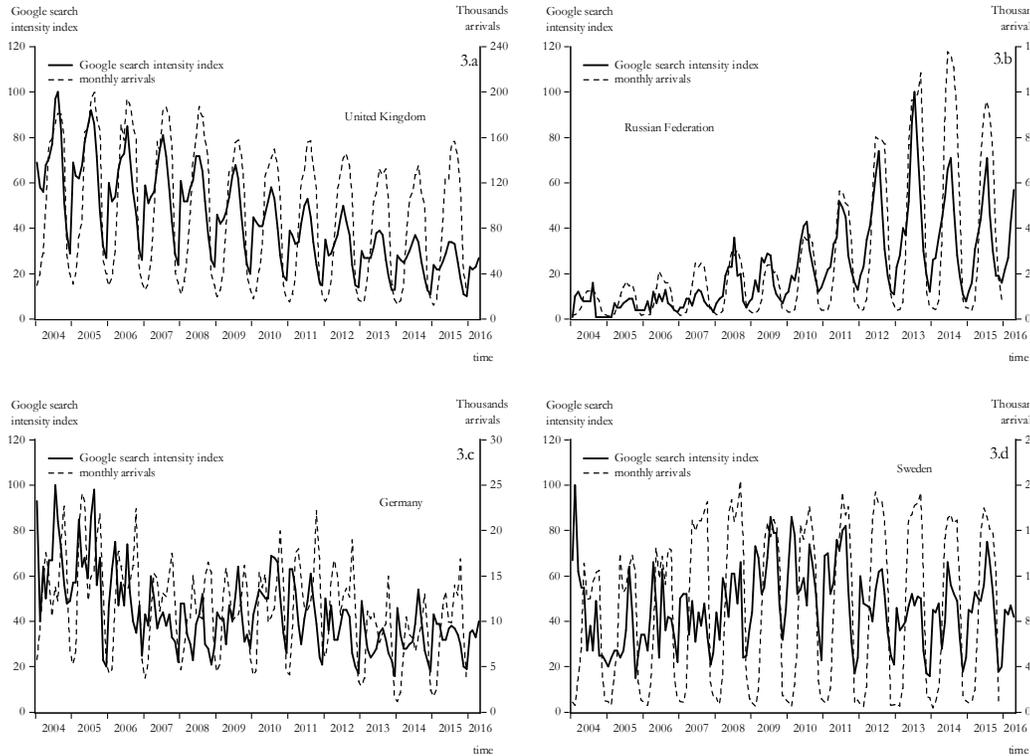


Figure 3. Google *SII* and tourists' arrivals per country.

In the case of Russia, the *platform bias* is also essential. Since Google is not the dominant search engine, we run the risk to misidentify the precise interest of the tourists and its evolution over time. To cross-check the validity of our selected keyword (туры кипр) from Google, we execute the same identification strategy by using a similar facility offered by Yandex (Relevant Phrases). The delivered keyword is туры на кипр, which is almost identical to the keyword identified by Google Trends (туры кипр). To assess the evolution of the two keywords across the two search engines over time, we take advantage of another feature offered by Yandex, which delivers the absolute number of searches for a keyword of interest. The common sample correlation coefficient between the *SII* of Google Trends (туры кипр) and the number of searches in Yandex (туры на кипр) is 0.97 (The absolute number of a search in Yandex is available, on a monthly basis, for the past two years). Hence, we may argue that the *SII* obtained from the Google, despite its' relatively small share on the market, reveals the true pattern over time.

To construct the aggregate uncorrected, for the two sources of bias, *SII*, we combine all the previously identified keywords (one for each country) to a single search. The conducted search is: hotel Cyprus + туры кипр + hotel zypern + cypern resor. The constructed index is presented in Fig. 4.

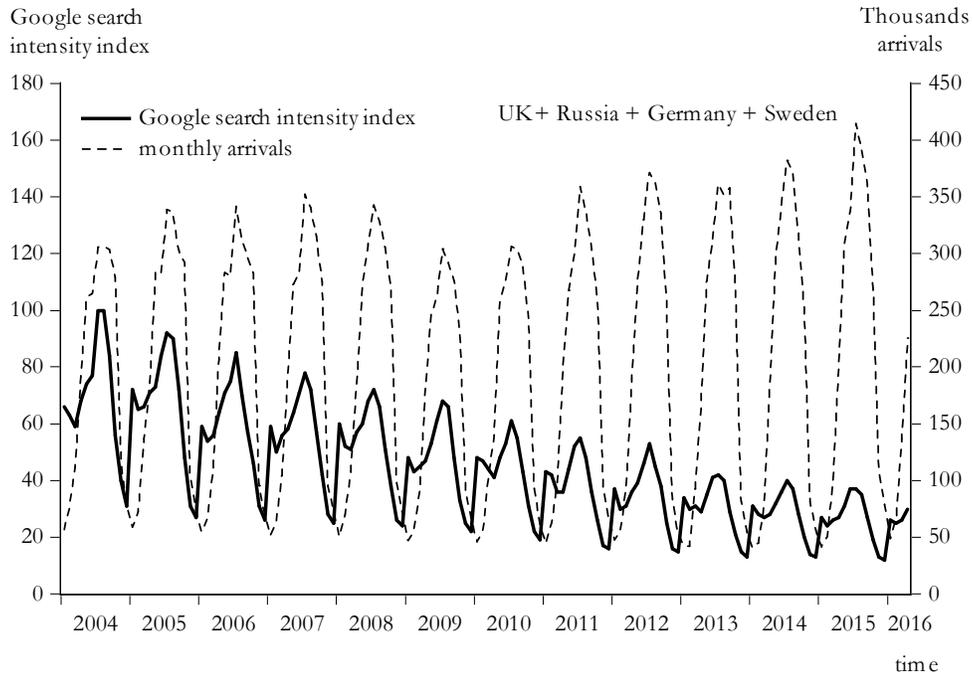


Figure 4. Aggregate Google *SII* and total arrivals of tourists.

4.2. Preliminary Econometric Analysis

The data of the arrivals and the web *SII* (see Figs. 3 and 4) show a clear seasonal variation. The well-known adverse effects of seasonality in statistical inference dictate a need for a seasonal adjustment procedure. In our case, we remove the deterministic seasonal parts of the series by implementing the TRAMO/SEATS approach as part of the X-13ARIMA-SEATS program. Fig. 5 illustrates the seasonally adjusted series.

In addition, the stationarity properties of all the de-seasonalized series are examined by conducting the Phillips and Perron (1988) test, with and without the presence of a deterministic linear trend (Table 1). Hence, we reject the null hypothesis of a unit root, at the 0.01 significance level, for the aggregate arrivals and the arrivals that originated from Germany and Sweden, while the opposite is true (failing to reject) for the arrivals that come from the U.K. and Russia. However, once we allow for the presence of a linear trend, the arrivals from the U.K. and Russia prove to be trend stationary. Similarly, the null hypothesis is rejected, at the 0.01 significance level, when the test is conducted to the Google *SII* of two countries, Germany and Sweden, while this is not the case for the remaining indices. The inference for the remaining indices is reversed in the presence of a linear trend. Overall, we may treat all the involved variables as stationary or trend

stationary. In the case where a variable is trend stationary, it is incorporated into our analysis after removing the linear time-trend.

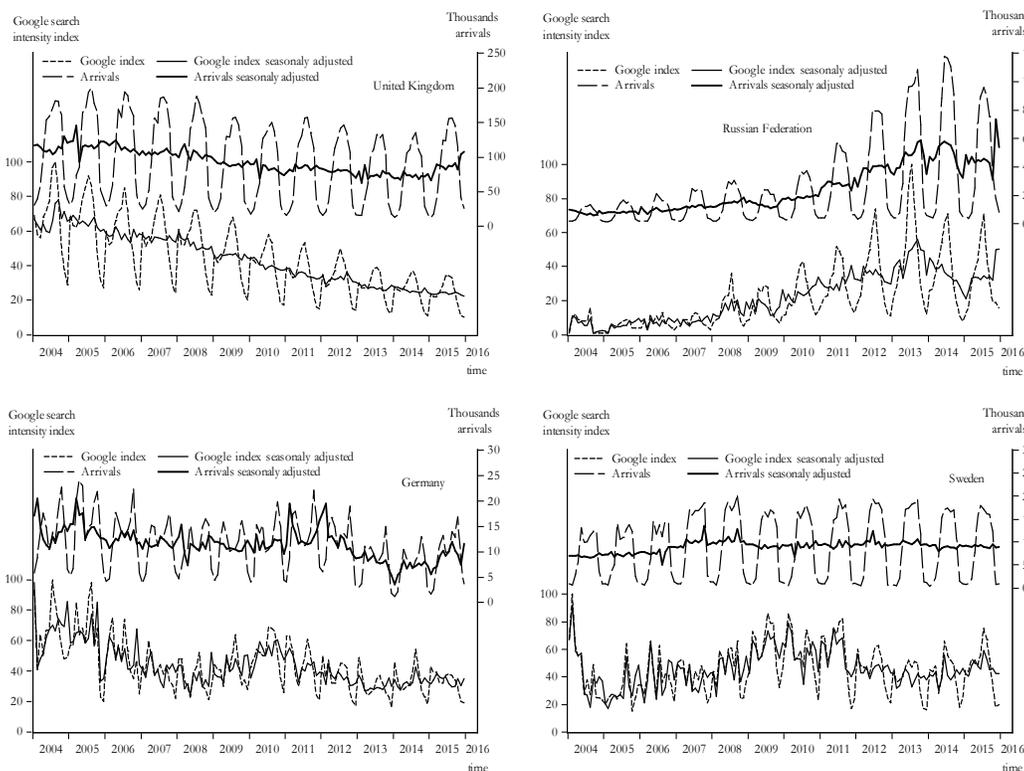


Figure 5. Seasonally adjusted series for the *SII* and the arrivals per country.

Table 1. Phillips-Perron unit-root tests for the de-seasonalized series.

Country	Arrivals		Inference	Google <i>SII</i>		Inference
	no trend	trend		no trend	trend	
UK	-1.75	-4.03***	I(0)/	-0.74	-3.83**	I(0)/
Russia	-0.44	-3.15*	I(0)/	-1.08	-3.17*	I(0)/
Germany	-3.58***	-4.68***	I(0)	-5.22***	-7.49***	I(0)
Sweden	-5.33***	-6.48***	I(0)	-6.21***	-6.28***	I(0)
Aggregate	-4.87***	-5.16***	I(0)	-0.58	-5.63***	I(0)/

Notes: the symbols * and *** denote the rejection of the null hypothesis at the 0.1 and 0.01 significance level, respectively. I(0) that implies that the series is stationary, while I(0)/ implies that the series is stationary under a linear time-trend. Finally, the bandwidth for the Phillips-Perron test was chosen based on the Newey-West selection procedure, while the spectral estimation method used is the Bartlett kernel.

5. Empirical Results

5.1 Predictive Power of the Web Search Intensity per Individual Country

To evaluate the predictive content of the constructed Google *SII* for each individual country towards the respective arrivals in Cyprus, we implement two alternative causality tests: the standard linear Granger non-causality test in the time domain and the B&C non-

causality test in the frequency domain. The B&C allows us to identify whether a verified causal relationship is short-run or long-run, and can reveal potential non-linear causal relationships. Additionally, working within the frequency domain could help to disclose causal relationships that may not be distinguishable in the time domain.

Table 2 shows that the hypothesis of no predictability running from the *SII* to the arrivals is consistently rejected for all countries of interest. In particular, predictability is verified at the 0.05 significance level for the U.K. and Sweden, while the same inference is drawn for Russia and Germany at the 0.01 significance level. Additionally, we fail to reject the hypothesis of no predictability that runs from the arrivals to the *SII*. The only exception is Sweden where bidirectional causality is established. Overall, our findings based on the standard Granger test suggest that arrivals in Cyprus from the four major source markets can be predicted by the respective *SII*.

Table 2. Standard Granger non-causality test results (per country).

Country	Google <i>SII</i> → Arrivals		Arrivals → Google <i>SII</i>	
	<i>F</i> -statistic	(lag length)	<i>F</i> -statistic	(lag length)
UK	3.67**	(3)	1.51	(3)
Russia	9.96***	(3)	1.52	(3)
Germany	4.54***	(4)	0.73	(4)
Sweden	2.31**	(5)	3.41***	(5)

Notes: the symbols ** and *** denote the rejection of the null hypothesis of non-causality at the 0.05 and 0.01 significance level, respectively. The numbers within the parentheses indicate the lag length of the underlying bivariate VAR specification. Finally, the arrow signifies the direction of causality.

The B&C test results for the U.K., in Fig. 6.a, show that the null hypothesis of no predictability running from *SII* to tourist arrivals is rejected at the 0.05 significance level, when $\omega \in [0, 1.24]$. This finding suggests that low and medium cyclical components of the *SII*, with wavelengths of more than five months, are those that contribute significantly in predicting arrivals. The opposite hypothesis is clearly rejected for the entire range of frequencies. The results for Russia are shown in Fig. 6.b. In particular, the predictability of the arrivals through *SII* is verified for the entire set of frequencies ($\omega \in [0, \pi]$). Again, the opposite hypothesis is rejected for the complete set of frequencies. For Germany (Fig. 6.c), predictability is not verified for the medium cyclical components but rather for the low and the high cyclical components of the series ($\omega \in [0, 0.75] \cup [1.88, \pi]$). Therefore, significant predictability is confirmed for wavelengths of less than 3.3 months and more than 8.4 months. Again, arrivals appear not to predict the *SII*. Finally, our findings for Sweden (see Fig. 6.d) show that only the high-frequency components of the *SII* series are

significant in predicting arrivals ($\omega \in [1.85, \pi]$). Hence, predictive power exists for wavelengths of less than 3.4 months. As was the case with the linear Granger non-causality test, we reject the non-predictability for the opposite hypothesis in high frequencies ($\omega \in [1.97, \pi]$), implying predictability for wavelengths of less than 3.2 months. In other words, for the case of Sweden, short-run bidirectional predictability is established. Overall, our findings from the B&C test are qualitatively similar to those of the linear Granger non-causality test.

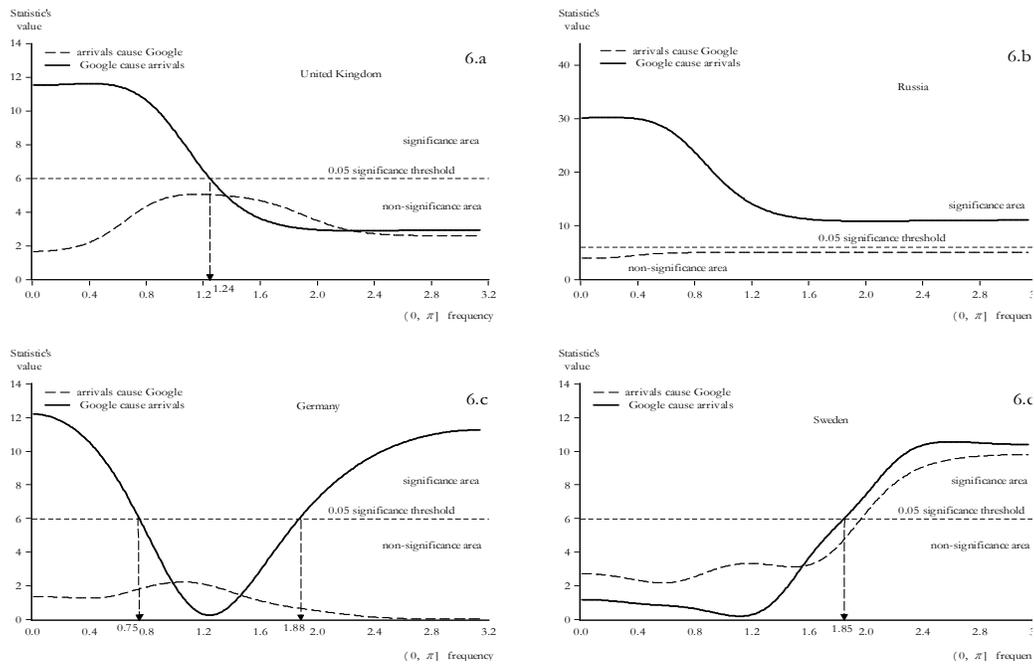


Figure 6. B&C Granger non-causality test per country

However, both non-causality tests are unable to reveal whether the variables of interest are connected in a positive or negative manner. The Cholesky defined accumulated impulse response functions of interest along with their associated ± 2 standard errors confidence bands are presented in Figs. 7a to 7d. For the case of the U.K., the accumulated response of tourist arrivals to one standard deviation shock in the *SII* for a 10-month period. Clearly, the response of the arrivals is constantly positive and significant for the entire period. Additionally, the impulse response analysis supports further our findings in the B&C test for the existence of causality that is long-run in nature. The impulse response analysis for Russia (see Fig. 7.b) and Germany (see Fig. 7.c) provides qualitatively similar inference to that of the U.K.. Hence, we observe a constantly positive and significant

response of the arrivals to one standard deviation shock in the *SII* for both countries. Finally, for the case of Sweden (see Fig. 7.d), the impulse response function is positive throughout the examined period, but it proves to be significant only in the first few months. This finding is consistent with the B&C test results which support causality only in the short-run. Overall, we may claim that the response of the arrivals in Cyprus to one standard deviation shock in the *SII* is positive and in harmony with the B&C test results.

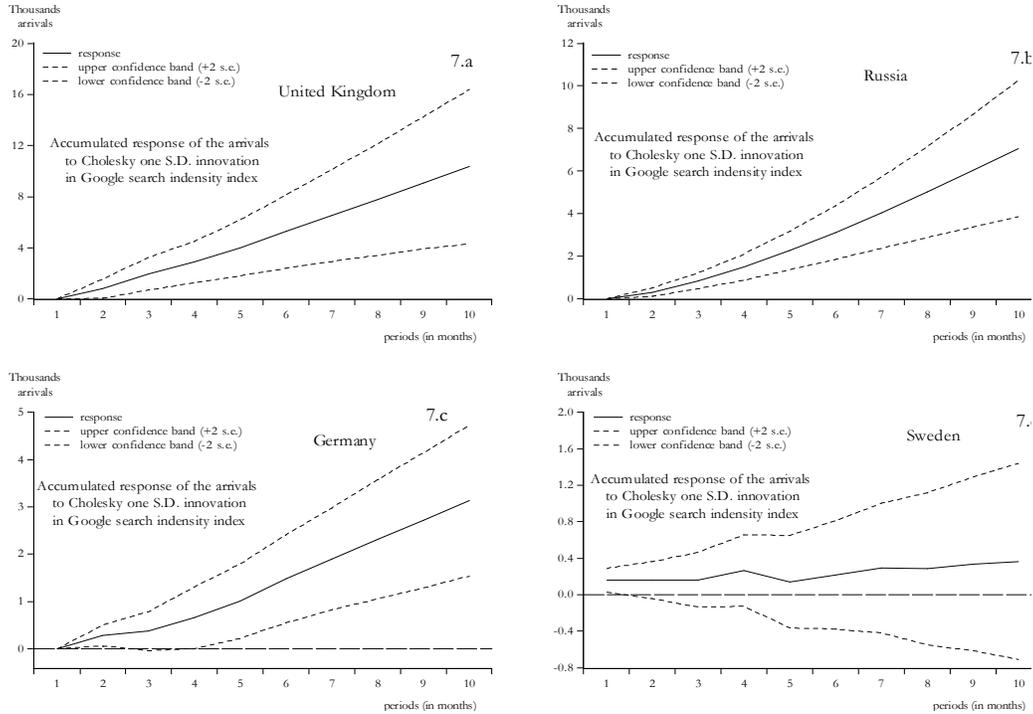


Figure 7. Impulse response functions per country

5.2 Predictive Power of Uncorrected Aggregate Web Search Intensity Index

Next, we perform the causality tests for the aggregate web *SII* concerning the total arrivals (see Fig. 4). After de-seasonalizing both series and de-trending the aggregate web *SII*³ (see the unit-root test results in Table 1), we conduct the standard Granger non-causality test (Table 3). We fail to reject the null hypothesis of no predictability that runs from the aggregate *SII* to the total arrivals (see 1st line in Table 3). Finally, the same inference holds for the opposite hypothesis.

³ To save space these results are not presented here. They are available upon request.

Table 3. Standard Granger non-causality test results (aggregate).

Country	Google <i>SII</i> → Arrivals		Arrivals → Google <i>SII</i>	
	<i>F</i> -statistic	(lag length)	<i>F</i> -statistic	(lag length)
Aggregate	1.37	(3)	0.03	(3)
Aggregate corrected	4.09***	(3)	1.07	(3)

Notes: the symbols ** and *** denote the rejection of the null hypothesis of non-causality at the 0.05 and 0.01 significance level, respectively. The numbers within the parentheses indicate the lag length of the underlying bivariate VAR specification. Finally, the arrow signifies the direction of causality.

The hypothesis of no predictability is also certified within the framework of the B&C test. In particular, the null hypothesis of no predictability running from the *SII* to tourist arrivals is not rejected, at the conventional levels of significance, for the entire set of frequencies ($\omega \in [0, \pi]$). Similarly, arrivals fail to predict *SII* in any significant manner the (see Fig. 8.a). Although the associated impulse responses are proved consistently positive, the relevant confidence bands include the zero value throughout the examined period. These findings are consistent with the B&C test results. Overall, while there is convincing evidence of predictability at a country level, this predictability vanishes once we use the aggregate data.

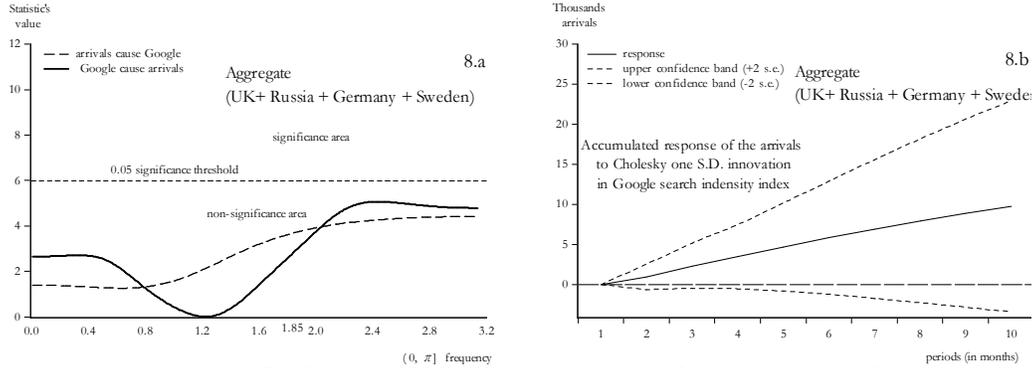


Figure 8. B&C test and impulse responses at the aggregate level.

5.3 Predictive Power of Corrected Aggregate Web Search Intensity Index

To reassess the predictive content of the aggregate index, we construct a corrected index which considers the *language bias* and *platform bias*. In particular, we follow two steps in constructing the aggregate corrected index. First, to overcome the *language bias*, instead of using the total number of tourist arrivals, we restrict our exercise only to the arrivals that correspond to the four major source markets (U.K., Russia, Germany and Sweden, which jointly pose almost 70% of the overall share in the arrivals). The constructed aggregate index reflects truly the web search intensity which is linked to the arrivals from these countries.

Second, to rectify our aggregate index from the *platform bias*, which is present in the case of Russia, we need to correct for the low market share of Google in the Russian Internet market. Instead of constructing a unified index by combining our four keywords (hotel Cyprus + туры кипр + hotel zypern + cypern resor), we extract four separate indices (one for each keyword) which are now compared jointly in terms of search volume (See Fig. 9). As Google has a low market share (approx. 25%, S_1) in the Russian internet market, naturally the *SII* that corresponds to Russia (see Fig. 9.a), underestimates the true search intensity.

At the same time, as Yandex dominates the Russian internet market (with a market share approx. 60%, S_2) and provided that the volume delivered from Yandex (for the keyword туры на кипр) correlates strongly to the index delivered from Google (for the keyword туры кипр), we may use the ratio of the respective market shares (S_2/S_1) as a market share correction factor. Once we multiply Google’s web *SII* that corresponds to Russia with the correction factor (S_2/S_1), then we can add the corrected index for Russia to the remaining three indices to form the aggregate corrected index. Consequently, the corrected aggregate index is expected to receive values above 100. This scale adjustment is attributed to the alternative scaling factor as well as to the introduced correction factor (these details are analytically discussed in the Appendix). For comparison purposes, the aggregate corrected *SII* along with the initial aggregate *SII*, both are illustrated in Fig. 9.b.

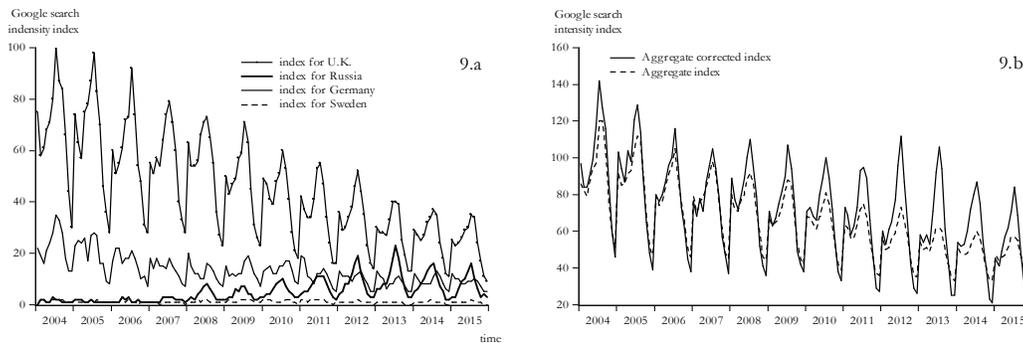


Figure 9. Search volume of the selected keywords and the aggregate corrected index.

Working within the same methodological framework, we can examine the predictive content of the corrected aggregate *SII* (See Fig. 9.b) with respect to the total arrivals from the four main source markets. Before testing for non-causality, we de-seasonalize the arrivals from the four major source markets and the corrected aggregate index, while we

de-trend only the latter. Starting from the standard linear non-causality test, we now fail to reject the hypothesis of no predictability that runs from the corrected aggregate index to the total arrivals from the four major source markets (see the 2nd line in Table 3) for all the conventional levels of significance. Regarding the opposite hypothesis, the testing results imply no predictability. The B&C test shows qualitatively analogous inference. The predictability running from the corrected aggregate index to the total arrivals from the four major source markets is verified at the 0.05 significance level, for wavelengths of more than 3.6 months ($\omega \in [0, 1.73]$) (See Fig. 10.a), while for the opposite hypothesis, there is no predictability at any frequency. Finally, the associated impulse response function is consistently positive with the confidence bands not to include the zero value (See Fig. 10.b).

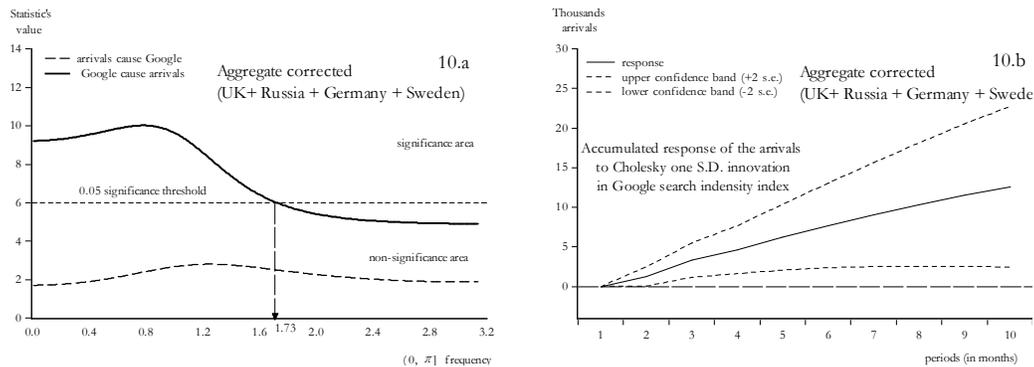


Figure 10. B&C test and impulse responses at the aggregate level (corrected SII).

6. Behavioral and Managerial Implications

Undeniably, search engines are among the most popular online planning sources for travelers (Google Travel Study, June 2014, Ipsos MediaCT⁴). The present study, focusing on search engines, offers valuable insights to the understanding of travelers' behavior when planning to visit a destination. We argue that since web activity reveals visitors' preferences in real-time, and encompasses even the effect of sudden one-off events, web search activity indices should always be included as an explanatory factor in tourism demand functions. Our study reveals a crucial methodological aspect when aggregate search intensity indexes are constructed, as we flag that aggregate indexes are subject to two sources of bias, the *language bias* and the *platform bias*. Therefore, we validate that when one predicts the

⁴ See: https://storage.googleapis.com/think/docs/2014-travelers-road-to-decision_research_studies.pdf

consumption of the tourist product based on aggregate web indices, bias correction approaches should be implemented. Overall, our findings may assist researchers to specify more complete tourism demand functions by including the early signals coming from the web and at the same time to construct more robust web search intensity indices. This way, the resulted demand functions will be better specified leading to more accurate forecasts.

The results also reveal interesting travel planning behavior of tourists to Cyprus. The decision-making process of travelers from different countries is not following a uniform pattern, as national cultures influence the search information behavior (Mc Cabe *et al.*, 2016; Gursoy and Umbreit, 2004). Specifically, the results show that the predictive content of the constructed Google *SII* (with respect to the arrivals) is dissimilar among the examined countries. According to Mc Cabe *et al.*, (2016), national cultures represent idiosyncratic features, which affect the search information behavior. Similar in nature are the findings of Gursoy and Umbreit (2004), who verify for a set of European countries that national cultures influence traveller's search behavior, resulting in this way to distinct consuming patterns.

Heterogeneous consuming patterns imply variation at the lead-time of decision-making. In particular, our findings reveal for the three major source markets (U.K., Russia, and Germany) that a large amount of tourists choose Cyprus as a destination at least a half year ahead of their arrival time. A characteristic example is the case of Germany. According to the Reise Monitor survey conducted by the ADAC Verlag,⁵ which investigates the holiday travel patterns of German tourists, 70% of the travellers intending to visit European destinations start planning their trip half a year ahead. Similarly, the percentage of tourists who plan their trip three months ahead until the last minute is approximately 20%. Such pattern does corroborate our empirical findings. The *SII* for tourists from Sweden has a predictive content for wavelengths of less than 3.4 months. This pattern can be attributed to the idiosyncratic features of those Swedish tourists whose booking practices may be heavily depended on travel agencies; therefore, personal search for additional information may take place only a few months prior their trip.

In more detail, the findings of the present study might be useful to governmental agencies, stakeholders of the sector, and specifically Destination Management Organizations (DMO's), when their purpose is to identify the upcoming future demand. The Cypriot government, with improved knowledge on the total magnitude of the arrivals, can assess more accurately the sectors' contribution to the economy. Therefore,

⁵ The survey is available at: <http://www.pot.gov.pl/component/rubberdoc/doc/1897/raw>

projections about the country's future growth path involve lower uncertainty (see Clerides and Adamou, 2010). Moreover, our results may help the policymakers of the Cyprus Tourism Organization to attend tourism exhibitions and to conduct advertising campaigns on the right timing for each source market. In other words, by knowing the proper time that every action needs to take place, policy makers could achieve cost savings and efficient allocation of resources. Furthermore, prediction of decreased arrivals from one destination can assist in adopting promotion actions in other promising markets. This way, reduced arrivals can be reversed by increasing last-minute bookings.

Overall, knowledge about the upcoming trends in the arrivals along with the unveiled behavioral patterns of the tourists from the major source markets may help the tourism sector to improve the quality of the provided services and allow potential investors to plan their projects (e.g. development of infrastructures) with greater certainty. The government, as well as all the stakeholders of the sector, would be more informed to allocate limited resources more effectively and to plan short and long-run promotion and investment strategies.

7. Conclusions

In recent years, Google Inc. provides data on the intensity of queries conducted on their search engine. This leads to an outbreak of scientific inquiries aiming to explain upcoming trends in various markets based on these data. As search engines constitute a leading tool in planning vacations, the digital traces can be exploited to improve predictions on the consumption of tourism products. Under this prism, we examine the predictive power of a relevant web *SII*, as captured by Google, on the total number of arrivals at a destination of interest. While existing studies emphasize at destinations that receive arrivals from countries with one major language and Google to be the dominant web search platform (Bangwayo-Skeete and Skeete, 2015), our work is the first that focuses on a destination with a multilingual set of source markets and with different dominant search platforms. As such, we introduce an approach to correct for the *language bias* and the *platform bias*, improving the predictive power of the constructed index.

We test our hypothesis by using monthly tourist arrival data (2004–2015) to Cyprus and by conducting two Granger non-causality tests, the standard linear Granger non-causality test and the B&C non-causality test. By introducing a simple way to select

appropriate keywords and working within the above framework, our findings show: a) country-specific *SII* (for U.K., Russia, Germany and Sweden) are highly significant in predicting the arrivals from the corresponding source markets, b) the initially constructed aggregate *SII*, without considering the *language bias* and the *platform bias*, proves inadequate to predict the total number of arrivals, and finally c) a corrected version of the aggregate *SII*, taking into account the two problems, predicts the arrivals in a significant and positive manner (U.K., Russia, Germany and Sweden). A natural extension of this work is the complementary verification of the predictive capacity that the web *SII* carries towards arrivals for other destinations, through the execution of competing pseudo-forecasting exercises with different econometric specifications.

Overall, our study validates the usage of *SII* as an important leading indicator for the upcoming arrivals at a destination but also reveals one very crucial methodological aspect. For destinations that accept arrivals from countries in different languages, the formation of a precise aggregate *SII* (intended to capture the entire web activity) is a challenging task and in several cases almost impossible to be constructed. Therefore, we argue that when it comes to predicting the consumption of the tourist product based on the *SII*, then it is preferable to execute forecasting at a disaggregated level. In other words, every major source market has to be investigated separately. Acting such, we use a richer set of information by taking into account each country's idiosyncratic characteristics. Clearly, we do not claim that approaches aiming to predict arrivals at an aggregate level have to be ostracized. Instead, we support that aggregate *SII* are exposed to two significant biases, and hence special handling is needed. In failing to account for these biases, misleading prediction inferences may be conducted.

Acknowledgements

We would like to thank the Editor in Chief Chris Ryan, and two anonymous reviewers of this Journal for their most useful comments and suggestions. Dr. Bing Pan would like to acknowledge the grant support from the National Natural Science Foundation of China (Grant # 41428101).

References

- Bangwayo-Skeete, P.F. and R.W. Skeete (2015). Can Google data improve the forecasting performance of tourist arrivals? Mixed-data sampling approach. *Tourism Management* **46**(1), 454-464.
- Breitung, J. and B. Candelon (2006). Testing for short- and long-run causality: A frequency-domain approach. *Journal of Econometrics* **132**(2), 363-378.
- Beracha, E. and M.B. Wintoki (2013). Forecasting residential real estate price changes from online search activity. *Journal of Real Estate Research* **35**(3), 283-312.
- Burger, C., M. Dohnal, M. Kathrada and R. Law (2001). A practitioners guide to time-series method for tourism demand forecasting - a case study of Durban, South Africa. *Tourism Management* **22**(4), 403-409.
- Chen, K.Y. and C.H. Wang (2007). Support vector regression with genetic algorithms in forecasting tourism demand. *Tourism Management* **28**(1), 215-226.
- Choi, H. and H. Varian (2009). Predicting the Present with Google Trends. Technical report, Google Inc.
- Choi, H. and H. Varian (2012). Predicting the Present with Google Trends. *Economic Record* **88**(281, S1), 2-9.
- Clerides, S. and A. Adamou (2010). Prospects and Limits of Tourism-Led Growth: The International Evidence. *Review of Economic Analysis* **2**(3), 287-303.
- Craigwell, R. and M. Winston (2008). Foreign Direct Investment and Tourism in SIDS: Evidence from Panel Causality Tests. *Tourism Analysis*, **13**(4), 427-432.
- Da, Z., J. Engelberg and P. Gao (2011). In search of attention. *Journal of Finance* **66**(5), 1461-1499.
- Das, J. and C. Dirienzo (2010). Tourism Competitiveness and Corruption: A Cross-country Analysis. *Tourism Economics* **16**(3), 477-492.
- Deng, J., B. King and T. Bauer (2002). Evaluating natural attractions for tourism. *Annals of tourism research* **29**(2), 422-438.
- Dergiades, T., C. Milas and T. Panagiotidis (2015). Tweets, Google trends, and sovereign spreads in the GIIPS. *Oxford Economic Papers* **67**(2), 406-432.
- Dragouni, M., G. Filis, K. Gavrilidis and D. Santamaria (2016). Sentiment, mood and outbound tourism demand. *Annals of Tourism Research* **60**, 80-96.
- Fesenmaier, D., Z. Xiang, B. Pan and R. Law (2011). A framework of search engine use for travel planning. *Journal of Travel Research* **50**(6), 587-601.
- Frechtling, D.C. (2001). *Forecasting Tourism Demand: Methods and Strategies*.: Butterworth-Heinemann. Oxford, UK.

- Gafer L.M. and A. Tchetchik (2017). The role of social ties and communication technologies in visiting friends' tourism: A GMM simultaneous equations approach. *Tourism Management* **61**(1), 343-353.
- Geurts, M.D. and I.B. Ibrahim (1975). Comparing the Box-Jenkins approach with the exponentially smoothed forecasting: model application to Hawaii tourists. *Journal of Marketing Research* **12**(2), 182-188.
- Goh, C., R. Law and H.M.K. Mok (2008). Analyzing and forecasting tourism demand: a rough sets approach. *Journal of Travel Research* **46**(3), 327-338.
- Granger, C.W.J. (1969). Investigating causal relations by econometric models and cross spectral methods. *Econometrica* **37**(3), 424-438.
- Gursoy, D. and T. Umbreit (2004). Tourist information search behavior: cross-cultural comparison of European Union member states. *International Journal of Hospitality Management* **23**(1), 55-70
- Halicioglu, F. (2010). An econometric analysis of the aggregate outbound tourism demand of Turkey. *Tourism Economics* **16**(1), 83-97.
- Joseph, K., M.B. Wintoki and Z. Zhang (2011). Forecasting abnormal stock returns and trading volume using investor sentiment: evidence from online search. *International Journal of Forecasting* **27**(4), 1116-1127.
- Kennedy, A.F. and K.M. Hauksson (2012). *Global Search Engine Marketing: Fine-Tuning Your International Search Engine Results*. Indianapolis: Que Publishing.
- Martin, C.A. and S.F. Witt (1989). Forecasting tourism demand: a comparison of the accuracy of several quantitative methods. *International Journal of Forecasting* **5**(1), 7-19.
- Martins, L.F., Y. Gan and A. Ferreira-Lopes (2017). An empirical analysis of the influence of macroeconomic determinants on World tourism demand. *Tourism Management* **61**(1), 248-260.
- Mc Cabe, S., C. Li and Z. Chen (2016). Time for Radical Reappraisal of Tourist Decision Making? Toward a New Conceptual Model. *Journal of Travel Research* **55**(1), 3-15.
- Meleddu, M and M. Pulina (2016). Evaluation of individuals' intention to pay a premium price for ecotourism: An exploratory study. *Journal of Behavioral and Experimental Economics* **65**(1), 67-78.
- Middleton, VTC, A. Fyall and M. Morgan (2009). *Marketing in travel and tourism*, 4th edition, Butterworth-Heinemann, Oxford, UK.
- Peng, B., H. Song and G. Crouch (2014). A meta-analysis of international tourism demand forecasting and implications for practice. *Tourism Management* **45**(1), 181-193.
- Phillips, P. and P. Perron (1988). Testing for a Unit Root in Time Series Regression. *Biometrika* **75**(2), 335-346.

- Poprawe, M. (2015). A panel data analysis of the effect of corruption on tourism. *Applied Economics* **23**(47), 2399-2412.
- Saha Ca, S. and G. Yap (2015). Corruption and Tourism: An Empirical Investigation in a Non-linear Framework. *International Journal of Tourism Research* **17**(3), 272-281.
- Sims, C.A. (1980). Macroeconomics and reality. *Econometrica* **48**(1), 1-48.
- Smith, G.P. (2012). Google Internet search activity and volatility prediction in the market for foreign currency. *Finance Research Letters* **9**(2), 103-10.
- Song, H., S.F. Witt and T.C. Jensen (2006). Forecasting international tourist flows to Macau. *Tourism Management* **27**(2), 214-224.
- Song, H., S.F. Witt and T.C. Jensen (2003). Tourism forecasting: accuracy of alternative econometric models. *International Journal of Forecasting* **19**(1), 123-141.
- Tang, S, E.A. Selvanathan and S. Selvanathan (2007). The relationship between foreign direct investment and tourism: empirical evidence from China. *Tourism Economics* **13** (1), 25-39.
- Turner L.W. and S.F. Witt (2001). Factors influencing demand for international tourism: tourism demand analysis using structural equation modelling, revisited. *Tourism Economics* **7**(1), 21-38.
- UNWTO (2012). Tourism in the green economy-background report. United Nations Environment Programme and World Tourism Organization, Madrid, Spain.
- Xiang, Z. and B. Pan (2011). Travel queries on cities in the United States: Implications for search engine marketing for tourist destinations. *Tourism Management* **32**(1), 88-97.
- Yang, Y., B. Pan and H. Song (2014). Predicting hotel demand using destination marketing organization's web traffic data. *Journal of Travel Research* **53**(4), 433-447.
- Yang, X., B. Pan, J. Evans and B. Lv (2015). Forecasting Chinese tourist volume with search engine data. *Tourism Management* **46**(1), 386-397.

Appendix

To construct the corrected aggregate intensity index, instead of conducting a joint search for the four keywords of interest (hotel Cyprus + туры кипр + hotel zypern + cypern resor), we perform a separate search by sequentially adding all the keywords of interest (see the *compare multiple search terms* in the help function of Google trends). Acting this way, we receive four separate series, which are directly comparable in terms of search volume (only one series receives the maximum value of 100). Having extracted the raw data for each search phrase, we adjust the series of interest with the market share correction factor, and then the four-separate series are added to form a single index.

Let's assume that we wish to compare four keywords. The search volume for each one of the queries, for the period of interest ($t=1,2,\dots,n$), can be denoted as: $V_{1,t}^q, V_{2,t}^q, V_{3,t}^q$ and $V_{4,t}^q$, respectively or more compactly as $V_{i,t}^q$ ($i=1,2,3,4$). Let now $V_{e,t}^q$ to represent, at time t , the entire volume of queries, then the first step of the normalization process that Google implements is to express the search volume of each query ($V_{i,t}^q$ with $i=1,2,3,4$) as a fraction of the entire search volume of queries ($V_{e,t}^q$), that is:

$$r_{1,t}, r_{2,t}, r_{3,t} \text{ and } r_{4,t} \text{ or } \frac{V_{i,t}^q}{V_{e,t}^q} = r_{i,t} \quad (i=1,2,3,4) \quad (\text{A.1})$$

Once the fractions have been estimated the four-normalized series can be constructed by multiplying each series with the scaling factor: $100/r^*$, where r^* is the maximum observed fraction among the fractions that come from the four-constructed series, that is:

$$\max_{r_{1,t}, r_{2,t}, r_{3,t}, r_{4,t} \text{ and } r_{i,t} \in \mathbb{R}^+} \{r_{1,t}, r_{2,t}, r_{3,t}, r_{4,t}\} = \{r^*\} \quad (\text{A.2})$$

The four normalized directly compared series can be denoted as: $S_{i,t}^n = (r_{i,t}/r^*)100$, with $i=1,2,3,4$. Once we have at our disposal the normalized series (this is the form that the Google Trends facility deliver's the series), we may now implicate the market share correction factor for the intensity index that corresponds to Russia, say $S_{4,t}^n = (r_{4,t}/r^*)100$. In particular, the volume adjusted series for Russia is now given by: $S_{4,t}^{n,va} = r_{4,t} (m/r^*)100$, where m is a scalar and represents the market share correction factor.

Given that the denominator is common, it comes that all four series can be added in order to form a unified, volume corrected, search intensity as follows:

$$S_t^f = \sum_{i=1}^3 S_{i,t}^n + S_{4,t}^{n,va} \text{ or } S_t^f = \frac{\sum_{i=1}^3 V_{i,t}^q + mV_{4,t}^q}{\frac{V_{e,t}^q}{r^*} + \frac{mV_{e,t}^q}{r^*}} 100 \quad (\text{A.3})$$

From A.3 it is obvious that it is possible to receive series that are scaled above 100. The difference of A.3 from the standard case, where the search of multiple keywords delivers a unique S_{II} with a maximum value of 100, lies on the fact that a) the scaling factor, r^* , is now different, and b) the market share correction factor is introduced. Given that both factors are simple scalars, the resulted series from the two alternative approaches are expected to illustrate almost identical evolution over time and therefore, a high degree of correlation. In other words, both approaches deliver qualitatively equivalent results.