# Questioni di Economia e Finanza

(Occasional Papers)

How can Big Data improve the quality of tourism statistics?
The Bank of Italy's experience in compiling the "travel" item
in the Balance of Payment

by Andrea Carboni, Costanza Catalano and Claudio Doria

# BANCA D'ITALIA

EUROSISTEMA

# Questioni di Economia e Finanza

(Occasional Papers)

How can Big Data improve the quality of tourism statistics?
The Bank of Italy's experience in compiling the "travel" item
in the Balance of Payment

by Andrea Carboni, Costanza Catalano and Claudio Doria

*The series* Occasional Papers *presents studies and documents on issues pertaining to the institutional tasks of the Bank of Italy and the Eurosystem. The* Occasional Papers *appear alongside the* Working Papers *series which are specifically aimed at providing original contributions to economic research.*

*The* Occasional Papers *include studies conducted within the Bank of Italy, sometimes in cooperation with the Eurosystem or other institutions. The views expressed in the studies are those of the authors and do not involve the responsibility of the institutions to which they belong.*

*The series is available online at www.bancaditalia.it .*

# HOW CAN BIG DATA IMPROVE THE QUALITY OF TOURISM STATISTICS? THE BANK OF ITALY'S EXPERIENCE IN COMPILING THE "TRAVEL" ITEM IN THE BALANCE OF PAYMENTS

by Andrea Carboni*, Costanza Catalano* and Claudio Doria*

## Abstract

In tourism statistics it is becoming more and more important to identify data sources that are more timely and cheaper than the traditional ones, such as surveys. In this paper, we investigate how mobile phone data (MPD), electronic payments data and internet search data (Google Trends) can improve the compilation of tourism statistics and the 'travel' item in the Balance of Payments (BoP). We find that MPD have the potential to improve the estimates of the number of international travelers and can be integrated with surveys, although a constant interaction with the data supplier is required to identify the phenomena to be captured. We highlight the limitations and issues in using electronic payment data for estimating expenditure in tourism statistics, and we propose a model for producing more timely preliminary estimates for BoP purposes. Finally, we point out that Google Trends data can be used to complement the sample estimates of international travelers and to improve the quality of provisional data.

## Contents

_____

* Bank of Italy, Statistical Analysis Directorate.

# 1 Introduction

The use of big data is rapidly spreading in several fields as economics and statistics. National Central Banks play a role in this growing exploitation, as in general, official statistics follow a pressing and strictly defined calendar.[1] Therefore, timely information, such as that coming from big data sources, is very attractive and potentially useful for compilers. Moreover, big data can be of great help as a supplementary source whenever information from traditional sources is difficult to obtain, time demanding and burdensome to acquire. In 2014, the UN Statistical Commission, recognizing the relevance of these new data sources, established the Global Working Group on Big Data for Official Statistics to promote the use of big data for compiling official statistics [29].

Against this scenario, the Bank of Italy has carried out an in-depth analysis to understand whether, and how, big data can enhance the data production of the BoP "travel" item. This paper illustrates the results of this research, focusing on mobile phone data, electronic payment data (credit/debit cards) and web research information (Google trends), and discusses how the traditional approach for the compilation of the "travel" item can be improved by their use.

The paper is structured as follows. The next section describes the Bank of Italy's traditional methodology for the compilation of the "travel" item. Sections 3, 4 and 5 respectively illustrate the three research paths – based on the use of mobile phone data, electronic payments data and Google Trends data - developed for improving and validating the compilation approach. The final section summarizes the main findings and conclusions.

# 2 The estimate of the "travel" item in the Italian Balance of Payments

The Balance of Payments (BoP) is a statement that summarizes the economic transactions between residents and non-residents during a specific period (as defined in [20]). The "travel" item of the BoP covers the monetary value of goods and services acquired in a country by non-resident travelers,[2] concerning visits to that country (see [20] §10.86), except the expenses for transport incurred to reach it,[3] which are instead recorded under the "passenger transport" item. The BoP compilation standards of the "travel" item require a breakdown by counterpart countries and by purpose of the visit (personal vs. business travels).

---

[1] E.g. EU Member States monthly disseminate provisional statistics of Balance of Payment at t+45 days.

[2] Only transactions between residents and non-residents are relevant. The definition of "residence" is economic and not administrative: the country of residence of an international traveler is where the center of her economic interest is located. To ease the exposition, we use for travelers the improper terms of "Italian" and "foreigner" when referring to their residence.

[3] E.g. flight tickets, international train tickets, tolls, etc. The expenses for transports within the visited economy are instead included in the "travel" item.

Countries can adopt different methodological approaches[4] for compiling this item, based on the relevance of tourism in their economy, the characteristics of the border points, the administrative controls on incoming and outgoing flows and, of course, the budget constraints.

Since 1996, the Bank of Italy has been collecting the relevant information (number of international travelers, expenditures, length of the trips) through a sample survey carried out at border points; the data collected are then integrated, when available, with administrative sources (e.g. the number of international travelers published by airports, ports authorities and railroads companies).

From a methodological point of view, the survey consists of two operations, carried out at each of the selected border points: counting and interviewing.

Counting aims at estimating, every month, the reference population, i.e. the total number of travelers entering or leaving Italy, broken down by country of residence or destination. In a selected interval of time, all the travelers crossing the border are counted and their residence is registered. Since having permanent counting operations on all the borders is not feasible, a grossing up algorithm is necessary for estimating the total amount of international travelers crossing the Italian borders during the reference period. Where administrative data are available, as for airports, they integrate the sample survey.

The second main operation consists in interviewing a sample of the travelers passing through each selected border point. The interviews primarily collect data on the expenditures and other relevant aspects for BoP purposes (e.g. reason of the trip, counterpart country), but also gather the information that allows a broader analysis of tourism related topics, such as the means of payment and the type of accommodation. The interviews are carried out at the end of the stay, which is when the memories of the traveler about the trip are the most recent and reliable possible, and all the expenses have already been determined so no guess by the traveler is necessary. Interviews are realized through questionnaires: each questionnaire refers to the group of people, if any, that shares the expenses of the trip with the interviewed traveler (e.g. a family).[5]

At each border point, interviewing and counting are carried out, as much as possible, at the same time, so that the characteristics of the interviewed traveler are coherent with those of the counted sample.[6] The information acquired with the questionnaires is then grossed up to the reference population, by taking into account the stratification variables listed in Table 1.

---

[4] In Europe, the main data sources are (frontier and households) sample surveys; some countries integrate this information with payment statistics as an additional source or for control purposes.

[5] For example, if the interviewed traveler shares the expenses with another traveler and they spend a total of 100 euros, we count two trips, each with a total expense of 50 euros.

[6] An interviewed traveler is not necessarily part of the counted ones. For example, on road borders it is very difficult to interview the passengers of the vehicles counted while passing through the frontier; it becomes possible only when there are checkpoints and collaboration with frontier authorities is arranged.

**Table 1.** Stratification variables and corresponding levels

| Variable | Levels |
|---|---|
| Direction | 2 (inbound, outbound) |
| Type of carrier | 4 (road, rail, airport, seaport) |
| Frontier point | 62 (22 roads, 4 rails, 25 airports, 11 seaports) |
| Day of data collection | number of days in the month (e.g. 31) |
| Time of the day | 3 (first shift, second shift, third shift) |

Annually about 100,000 interviews and 1,000,000 counting operations are realized in more than 60 frontier points. This size ensures that the sampling error of the total international travel expenditure estimates is small and the statistics for the main partner countries are reliable.

## 3 The mobile phone data experiment

Mobile Phone Data (MPD) are one of the most promising big data sources for the study of many social and economic phenomena and behaviors. Several pilot studies in the literature analyze the potentiality of MPD, e.g. for computing the population of an area [14], for estimating the population density [26], for traffic statistics [21], for transport and urban planning [24] and for travel statistics [1,2,3]. In this regard, the contributions of the Estonia Central Bank [17,18] and the Banque de France [22] have also to be mentioned.

While MPD data can provide a great amount of information (number of international travelers, a proxy of the country of residence, locations visited, length of stay, etc.) they say nothing about the expenses, the main variable to be estimated in the BoP "travel" item. MPD can thus only be considered as a complementary source of information, useful to estimate the dimension and some characteristics of the reference population.

In 2018, the Bank of Italy started a test phase to integrate MPD into the international frontier survey, to gradually replace the counting operations. Counting is a costly and demanding activity, and this is particularly true for road borders, given the high number of this type of frontier points in Italy and the scarcity of administrative data, and for seaports, due to restricted access zones as the ones reserved to cruise ships. These problems might affect the quality of the grossing-up factors and hence the estimated values.

MPD may represent an alternative, efficient and less costly data source to count travelers crossing the frontiers. The arrival of a foreign traveler at the Italian border is signaled by the connection of a mobile phone, with a SIM card[7] issued by a non-resident phone operator, to the cells controlled by an Italian phone-operator. Likewise, the disappearance for some time of the signal of an Italian SIM card near the border would indicate that this traveler has gone abroad.

---

[7] Subscriber Identity Module.

The Bank of Italy collaborated with one of the major Italian Mobile Network Operator[8] (MNO) to develop an algorithm for the estimate of travelers' inflows and outflows through each border point by exploiting the MPD. These data are not "ready to use" for BoP purposes and close, constant cooperation between the Bank of Italy and the MNO has been necessary to define the best metrics to elaborate the raw data and achieve measures compatible with the BoP standards. For example, it is necessary to define the minimum docking time of a foreign SIM card to a cell located in Italy for it to be considered an international traveler present in Italy. This problem is very relevant near road borders due to handover effects.

Since each frontier point has specific features that should be incorporated into the final algorithm, a test phase was developed for two important Italian border points: the main airport of Rome (Fiumicino), which is the largest in Italy in terms of international traffic, and the highway frontier point of Tarvisio, one of the most relevant in the North-East of Italy.

For the Fiumicino airport, the traditional survey is supported by data provided by the company that manages the airport, Aeroporti Di Roma (ADR). This source is used for correcting, using calibrated estimators (see [15]), the estimate of the total international flows derived from the counting operations, although it does not provide information on the residence of the passengers. Tables 2 and 3 compare the MPD, the ADR[9] statistics and the Bank of Italy's official statistics (BI) on the Fiumicino airport for the period August 2018 - June 2019.

**Table 2.** Fiumicino airport: comparison between MPD and ADR statistics on the number of international passengers.

|  | MPD[1] | ADR[2] | %[3] |
|---|---|---|---|
| Aug-18 | 1,802,051 | 1,679,511 | 7.3 |
| Sep-18 | 1,723,145 | 1,521,956 | 13.2 |
| Oct-18 | 1,590,179 | 1,437,316 | 10.6 |
| Nov-18 | 1,220,903 | 1,083,621 | 12.7 |
| Dec-18 | 1,045,675 | 1,066,898 | -2.0 |
| Jan-19 | 1,113,629 | 989.903 | 12.5 |
| Total | 8,495,582 | 7,779,205 | 9.2 |

[1] Estimates based on mobile phone data;
[2] Aeroporti di Roma administrative data on passenger transits at Fiumicino airport;
[3] Variation of MPD over ADR statistics.

The MPD and the ADR totals are broadly aligned, with the former almost always larger than the latter. As for the breakdown of residents/non-residents,

---

[8] 31% of the market share in 2018.
[9] Differences between the official data and ADR statistics are due to minor adjustment in the Bank of Italy estimation process.

**Table 3.** Fiumicino airport: comparison between MPD and Bank of Italy statistics on the number of international passengers.

| | Total | | | Italians | | | Foreigners | | |
|---|---|---|---|---|---|---|---|---|---|
| | BI[1] | MPD[2] | %[3] | BI[1] | MPD[2] | %[3] | BI[1] | MPD[2] | %[3] |
| Aug-18 | 1,717,076 | 1,802,051 | 4.9 | 640,288 | 621,419 | -2.9 | 1,076,788 | 1,180,632 | 9.6 |
| Sep-18 | 1,574,571 | 1,723,145 | 9.4 | 446,884 | 516,638 | 15.6 | 1,127,687 | 1,206,507 | 7.0 |
| Oct-18 | 1,380,639 | 1,590,179 | 15.2 | 423,402 | 449,204 | 6.1 | 957,237 | 1,140,975 | 19.2 |
| Nov-18 | 1,053,956 | 1,220,903 | 15.8 | 392,909 | 466,087 | 18.6 | 661,047 | 754,816 | 14.2 |
| Dec-18 | 1,037,503 | 1,045,675 | 0.8 | 506,530 | 417,820 | -17.5 | 530,973 | 627,855 | 18.2 |
| Jan-19 | 831,120 | 1,113,629 | 34.0 | 344,529 | 457,947 | 32.9 | 486,591 | 655,682 | 34.8 |
| Total | 7,594,865 | 8,495,582 | 11.9 | 2,754,542 | 2,929,115 | 6.3 | 4,840,323 | 5,566,467 | 15.0 |

[1] Bank of Italy official statistics;
[2] Estimates based on mobile phone data;
[3] Variation of MPD over BI statistics.

the number of Italian travelers estimated by the MPD is in line with the one estimated by the Bank of Italy, while for the number of foreign travelers, MPD are always greater than the Bank of Italy official statistics.

The estimate of the number of international travelers crossing the Tarvisio border only relies on counting operations, due to the lack of complementary administrative sources. Table 4 compares the Bank of Italy's and the MPD statistics in this road border point in the same period: the differences are very large, and they are wider for Italian travelers than for foreigners.

**Table 4.** Tarvisio: comparison between MPD and Bank of Italy statistics on the number of international passengers.

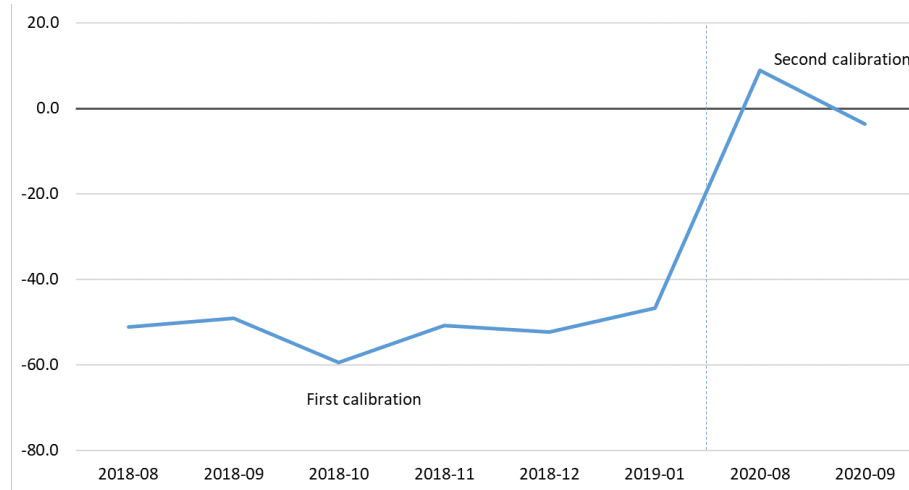| | Total | | | Italians | | Foreigners | |
|---|---|---|---|---|---|---|---|
| | BI[1] | MPD[2] | %[3] | BI[1] | MPD[2] | BI[1] | MPD[2] |
| Aug-18 | 2,005,595 | 980,066 | -51.1 | 662,710 | 115,841 | 1,342,885 | 864,225 |
| Sep-18 | 1,544,727 | 785,843 | -49.1 | 408,999 | 78,895 | 1,135,728 | 706,948 |
| Oct-18 | 1,026,265 | 416,988 | -59.4 | 261,135 | 64,948 | 765,130 | 352,040 |
| Nov-18 | 691,340 | 340,325 | -50.8 | 212,436 | 80,532 | 478,904 | 259,793 |
| Dec-18 | 600,309 | 285,953 | -52.4 | 272,065 | 103,197 | 328,244 | 182,756 |
| Jan-19 | 686,784 | 366,037 | -46.7 | 219,408 | 128,927 | 467,376 | 237,111 |
| Total | 6,555,020 | 3,175,212 | -51.6 | 2,036,753 | 572,340 | 4,518,267 | 2,602,873 |

[1] Bank of Italy official statistics;
[2] Estimates based on mobile phone data;
[3] Variation of MPD over BI statistics.

Further interactions with the mobile network operator led to a shortening (from four hours to 30 minutes) of the minimum time a foreign/Italian SIM card has to spend on the national/foreign territory to be considered an international traveler, and thus to a recalibration of the algorithm. This resulted in a new test,

only involving the months of August and September 2020: the result showed a good convergence on the total between the two data sources, with a great improvement compared to the first release (Figure 1). On the other hand, the distribution between resident and non-resident travelers was still quite different, suggesting the need to continue investigating the causes underneath different estimates.



**Fig. 1.** % differences between MPD and BI statistics on total travelers' flows.

## 4 The payment statistics analysis

Similarly to mobile phone data, electronic payment data are a promising source for the measurement and study of social and economic phenomena, including the production of statistics on national and international expenditure [12,13]. Recently, they have started to be used for tourism statistics [23], in particular by international institutes and national central banks such as the Banco de Portugal [10], the Banque de France [22] and the Central Bank of Armenia [32]. Moreover, the European Central Bank recently approved a regulation[10] on payment statistics to gather data that the Eurozone countries could use for the compilation of their external statistics [16].

Electronic payment data are attractive because of their timeliness, relative ease in collection and processing and moderate costs; moreover, their availability is not subject to high-impact perturbative phenomena like the Covid-19 pandemic. The steady increase of the share of electronic payments on total expenditure [4] will keep strengthening the informative power of this source, although

---

[10] Regulation ECB/2020/59.

all the other possible means of payment such as transactions made by cash, bank transfers, etc. have to be estimated with other sources.

Against this background, the Bank of Italy conducted an explorative analysis to assess if and to what extent electronic payment data can contribute to the production of the "travel" item of the BoP and/or can be used for checking the consistency with the tourism statistics.

For this purpose, two databases were considered, provided by one of the main paytech companies operating in Italy, with data spanning from May 2014 to August 2021. The market share of this company was unknown, making impossible the grossing-up of raw data. One database contains all the electronic payments made by credit and debit cards on POS[11] (physical database), while the other includes online (e-commerce) transactions. Both databases are divided into acquiring and issuing: the first one contains the transactions made by foreigners' cards on Italian POS and websites (potentially contributing to estimating the foreigners' expenditures in Italy), while the latter includes the transactions on foreign POS and websites made by Italian cards (potentially contributing to estimate the Italians' expenditures abroad).

Every record of the databases is made up of five variables: the date (day-month-year) of the payment, the Merchant Category Code (MCC) identifying the type of purchase[12], the nationality of the bank emitting the payment card, the country of the POS/website where the payment has been made and the amount of the transactions in euro.[13]

There are major limitations in electronic payments data for the compilation of official statistics on travel:

1. the nationality of the bank issuing the card is just a proxy of the residence of the traveler;
2. confidentiality issues allow the use of only aggregated data, which may increase the difficulties in discerning the transactions that are related to tourism from the ones that are not;
3. there is no information on the reason for the trip (business/personal), which is a mandatory BoP requirement;
4. there are difficulties in registering and correctly classifying the Digital International Platforms (DIP) transactions, in terms of misallocation issues for the counterpart country and of failure in recording some transactions. Three main examples help understanding the matter:
   (a) the payment of a stay in Paris made by an Italian tourist on the Booking.com platform[14] is recorded as a transaction from Italy to The Netherlands and not to France, as it should be recorded in the BoP;
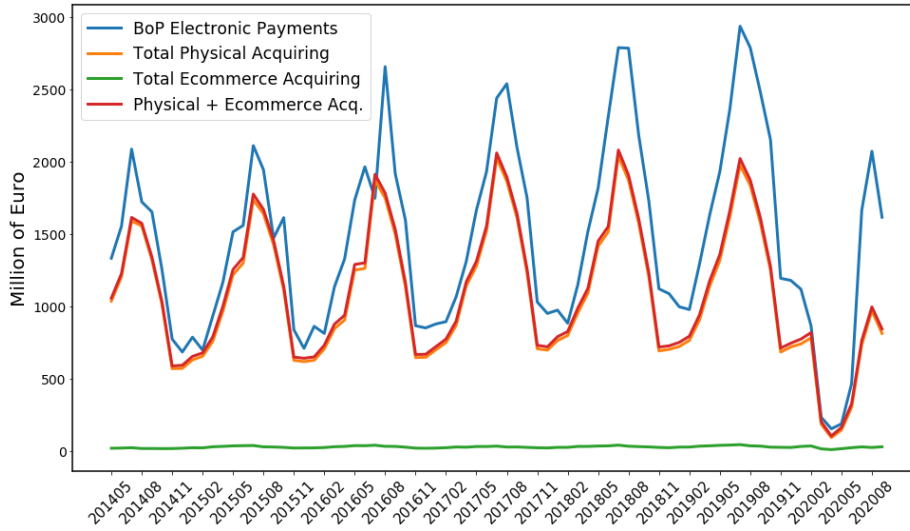
---

[11] Point Of Sale.

[12] The available MCC are the following: clothing, hotels and restaurants, groceries, home, cash advance, work, retail, services, mobile web, travels and transports.

[13] The amount is the aggregation of all the transactions sharing the same values of the first four variables.

[14] Whose legal headquarter is in The Netherlands.

(b) the payment on Airbnb[15] of accommodation in Rome by a French traveler is recorded as a transaction from France to Ireland, thus not appearing in our database, although it should be recorded in the BoP[16];

(c) the payment on Airbnb of accommodation in Rome by an Italian traveler is recorded as a transaction from Italy to Ireland, although it refers to a domestic trip and thus should not be recorded in the BoP.

Figure 2 compares the official BoP data on foreign travelers' expenditure in Italy (only by means of electronic cards) with the grand totals of, respectively, the electronic payments recorded in the physical acquiring database, in the e-commerce acquiring database, and the sum of the two. The level of the e-commerce transactions is much lower than in the other time series. Indeed, it does not cover the transactions of item 4.(b): the large use of these platforms, which mostly have foreign headquarters, can explain its negligible level.
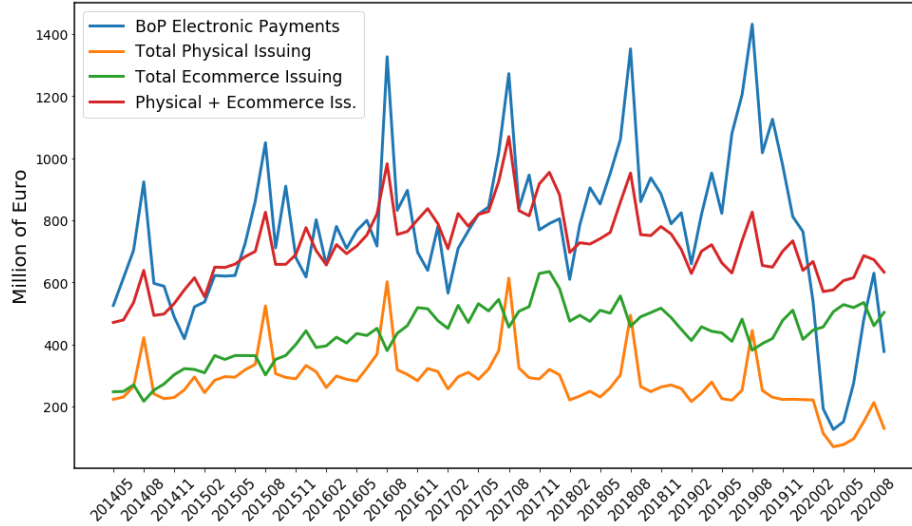


**Fig. 2.** Foreigner travelers' expenditure in Italy.

Figure 3 compares the official BoP data on Italian travelers' expenditure abroad (only by means of electronic cards) with the grand totals of, respectively, the electronic payments from the physical issuing database, the e-commerce issuing database and the sum of the two. The e-commerce time series shows higher levels than for the acquiring side, although it does not show the typical seasonality of the tourism phenomena, which has peaks in the summer. This is probably

---

[15] Whose legal headquarter is in Ireland.

[16] The digital platform can carry out a further transaction with an Italian counterpart, but not necessarily using a credit card.

because such a database contains large amounts of online transactions that are not connected to tourism, such as the purchases of goods on Amazon, Ebay, etc. The low granularity of the database does not always allow to distinguish them, as some categories contain both transactions that are related to travels and transactions that are not.



**Fig. 3.** Italian travelers' expenditure abroad.

In an attempt to push further the exploratory analysis, the following models have been tested: Ridge and Lasso[17] models, regression trees and boosted regression trees. Due to the relatively short length of the whole payment data time series,[18] the years 2015-2017 have been used as the training set, the year 2018 as the validation set[19] and the year 2019 as the test set, in order to verify the model out-of-sample performance by using the Mean Square Error (MSE) index.[20] Moreover, the Covid-19 years 2020-2021[21] were used as supplementary test set for verifying the robustness of the model to external shocks. In each model, the dependent variable is the total BoP travel "item", while the independent variables are all the physical MCC data (at lag 0) plus all the e-commerce

---

[17] With and without enforcing positive coefficients.

[18] 60 months from January 2015 to December 2019, plus 20 months (January 2020-August 2021) during the Covid-19 pandemic.

[19] It was used to select the best parameter $\lambda$ for the Lasso and Ridge models by minimizing the Mean square error on it.

[20] For the regression tree models, both 2018 and 2019 are used as test sets, fixing the maximum high of the tree to 4.

[21] We will call it the *Covid set*.

MCC data for all lags from 0 to -4, as on-line purchases can be made in advance with respect to the actual trip.

Table 5 reports the performance of such models in terms of the MSE index[22] on the validation, test and Covid set.

**Table 5.** Model performance in predicting the BoP "travel" item on different sets.

|  | Acquiring | | | Issuing | | |
|---|---|---|---|---|---|---|
|  | MSE val | MSE test | MSE covid | MSE val | MSE test | MSE covid |
| Ridge | 0.09 | 0.76 | 3.82 | 0.13 | 0.39 | 0.53 |
| Lasso | 0.04 | 0.29 | 1.74 | 0.08 | 0.57 | 0.15 |
| Lasso positive coeff. | 0.04 | 0.30 | 1.79 | 0.15 | 1.29 | 0.16 |
| Decision tree | 0.20 | 0.28 | 1.15 | 0.70 | 1.93 | 2.33 |
| Boosted tree | 0.10 | 0.32 | 1.23 | 0.48 | 2.04 | 2.60 |

val = validation set; test = test set; covid = Covid years 2020-2021;
Acquiring: foreigner travelers' expenditures in Italy;
Issuing: Italian travelers' expenditures abroad.


On the acquiring side, almost all the models show a quite good performance; in particular, the Lasso models and the regression tree obtain the smallest MSE on the test set and perform quite well in the Covid one. On the issuing side, the performance of each model is worse than in the acquiring case, as expected. The best results are obtained by the Lasso models, which have an unexpectedly good performance in the Covid set.

Figure 4 shows the plots of the forecasts in the test and Covid sets compared with the official BoP figures, for selected models. The graphs confirm what was already pointed out: in forecasting foreigners' travel expenses in Italy, we obtain a good performance on the test set, while the forecast behaves quite poorly in the subsequent years affected by the Covid pandemic. On the other hand, the forecast of the Italians' travel expenses abroad is worse in the test set, but surprisingly good in the Covid one, as the trend is fully captured.
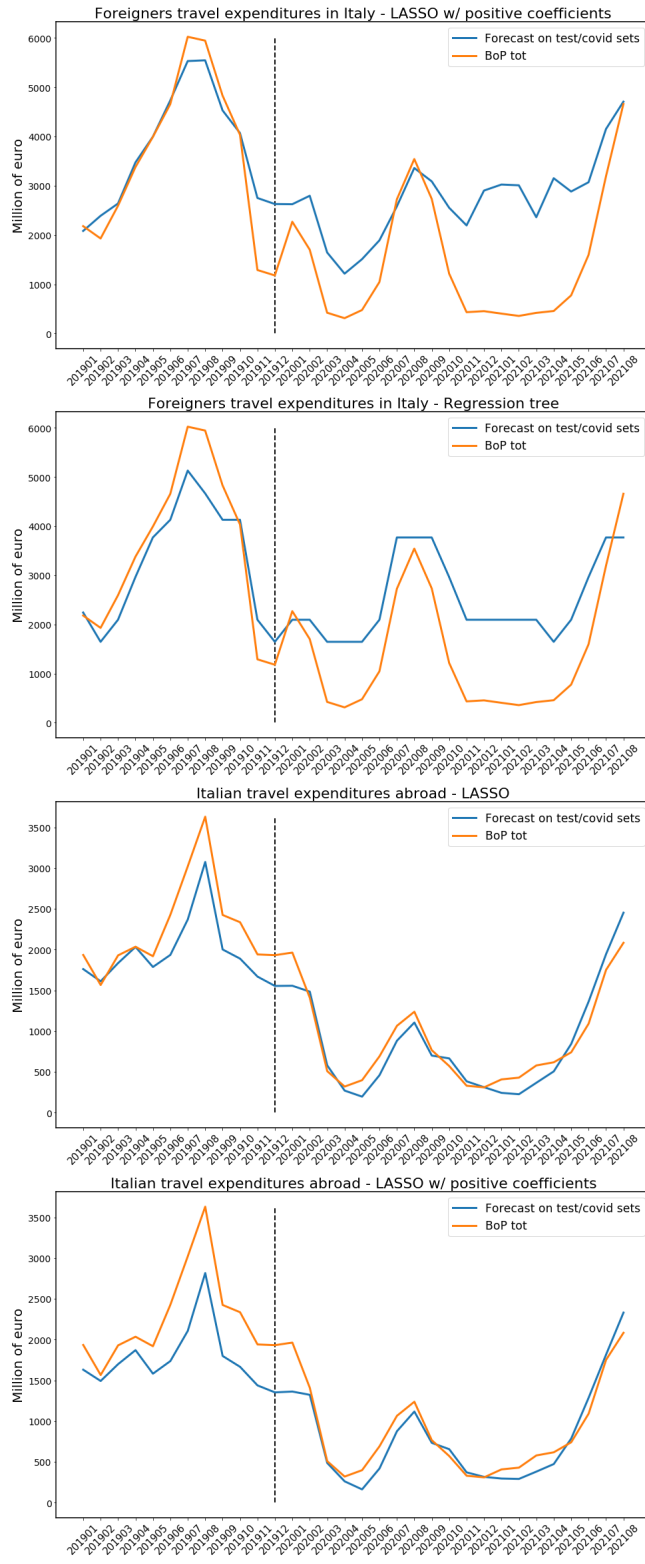

## 5   The Google Trends experiment

The third experimentation relies on the use of Google Trends as a complementary source to estimate the "travel" item in the BoP compilation process, in particular to assess the number of international travelers in Italy.

Google Trends ($GT$) is a website provided by Google [19] that reports the popularity of search queries in the Google search engine over time and across various regions of the world. The popularity of a given query is measured by an index between 0 and 100 (the maximum frequency). The data are collected and aggregated continuously on a daily, weekly or monthly basis. One can visualize

---

[22] We remind that the smaller the MSE is, the better the model performance is.

**Fig. 4.** Forecasts on test set and Covid set for selected models. Test set on the left of the black dashed vertical line, Covid set on the right.

the popularity of the selected query by specifying the state or region, the category they belong to,[23] and the time frame of interest. Timeliness is one of the main advantages of this website, as data are updated almost in real time.

To understand if this kind of data can be usefully employed, a specific predictive exercise was developed to forecast the number of foreigner travelers visiting Italy in the period from January 2006 to May 2019. In particular, we considered the tourist flows from the most important counterpart countries in terms of arrivals, namely France, Germany, United Kingdom, United States and Spain.

The $GT$ variable was defined by considering the frequency of the search queries performed in the aforementioned countries containing the word "Italy" in the "Travel" category . For each of these selected countries, a seasonal AR(1) process was used for modeling the number of travelers $N_{c,t}$ arrived in Italy during the month $t$ from country $c$ according to the Bank of Italy's tourism survey, where the $l$-period lagged Google Trends index $GT_{c,t-l}$ is included as an exogenous regressor:

$$N_{c,t} = \phi_0 + \phi_1 N_{c,t-1} + \phi_{12} N_{c,t-12} + \beta GT_{c,t-l} + \varepsilon_{c,t}. \tag{1}$$

The most suitable lag of the $GT$ index is chosen by minimizing errors of the out-of-sample forecasting performance, measured in terms of MSE reduction.

Figure 5 shows how the ratio between the MSE of specification (1) and the one obtained using the model without $GT_{c,t-l}$ ($\beta = 0$) depends on the lag for the different countries considered.

The contemporaneous variable $GT_{c,t}$ ($l = 0$) is the best predictor for Germany and Spain, while the variable at lag $l = 4$ and $l = 6$ minimizes the MSE for UK and US respectively. These last results seem only partially reasonable: US travelers may have to organize their trips towards Italy more in advance than German and Spanish travelers. Moreover, Google Trends may classify in the "Travel" category web searches performed by tourists during their travel, which should increase the weight for lag $l = 0$ in the model. Less clear is the situation for the UK, where there is not an intuitive explanation for a better performance of lag $l = 4$ in comparison to smaller lags.

The months between September 2012 and May 2019 have been considered to compare the observed value and the one-step-ahead forecasts, with an expanding windows approach.[24] In all cases, except for France where the coefficient $\beta$ is not statistically different from zero, the $GT$ index increased the performance of the predictive model.

The estimates of specification (1) for each country involved in the exercise are shown in Table 6: the $GT$ index is always highly significant and the model indicates a good in-sample fit, measured by a high value of the $R^2$.

---

[23] More than 20 categories of search, which helps avoid multiple meanings for the chosen query.

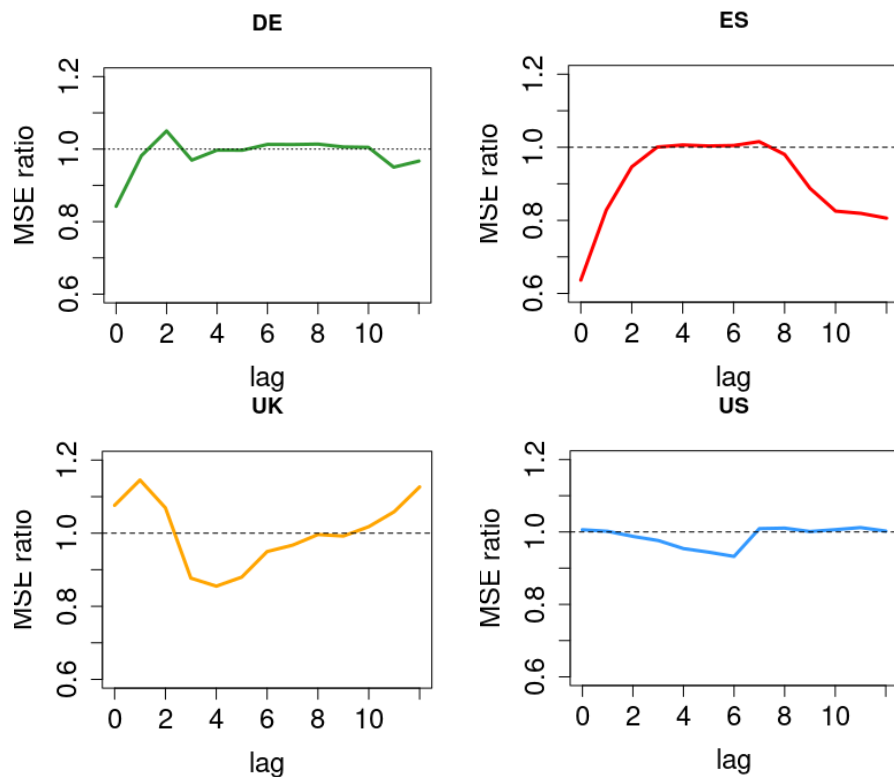[24] Adding an observation at each step.

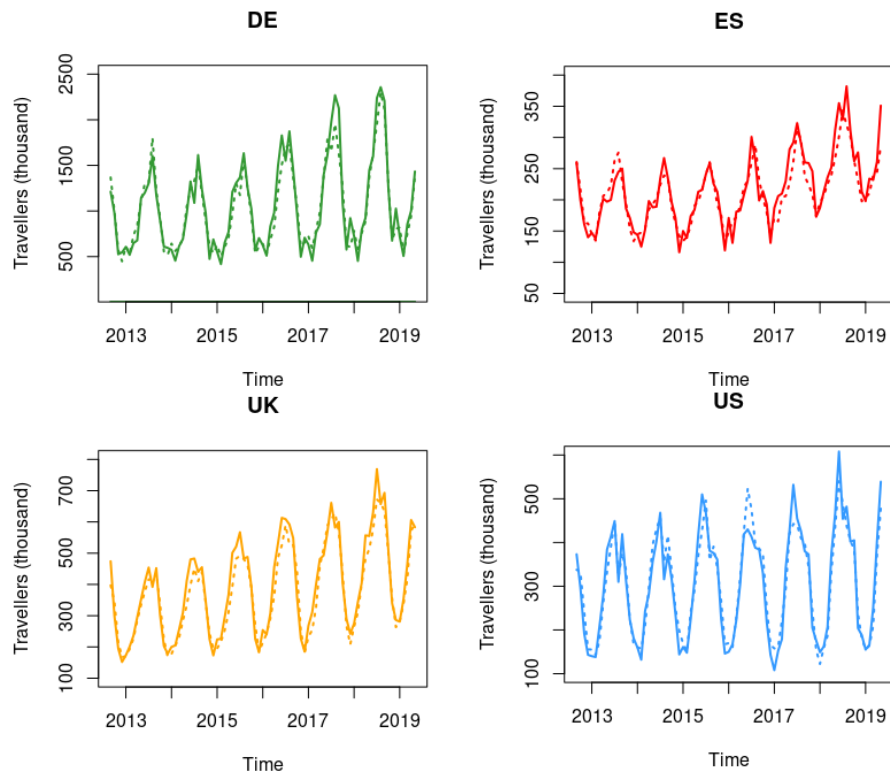**Fig. 5.** Out of sample normalized MSE for different lags of the $GT$ index.

**Table 6.** Estimates of the model for different countries.

|          | DE           | ES         | UK         | US         |
|----------|--------------|------------|------------|------------|
| $N_{t-1}$   | 0.20***      | 0.34***    | 0.33***    | 0.13***    |
|          | (0.04)       | (0.05)     | (0.05)     | (0.04)     |
| $N_{t-12}$  | 0.73***      | 0.46***    | 0.76***    | 0.89***    |
|          | (0.05)       | (0.05)     | (0.05)     | (0.04)     |
| $GT_{t-l}$  | 6.01***      | 2.11***    | -1.50***   | -0.91***   |
|          | (43.34)      | (12.07)    | (15.52)    | (13.93)    |
| Const    | -182.49***   | -22.78*    | 29.87*     | 38.84***   |
|          | (43.34)      | (12.07)    | (15.52)    | (13.93)    |
| $R^2$    | 0.92         | 0.77       | 0.89       | 0.92       |

$*p < 0.10, ** p < 0.05. *** p < 0.01$

However, each time series presents a strong auto-regressive component[25] and the marginal contribution of the $GT$ index is significant only for Spain.[26] The negative sign of the $GT$ coefficient in the UK and US regressions means that the variable is not robust enough for these two countries, confirming the doubts in the interpretation of the optimal lag.

To examine the predictive performance of the model Figure 6 compares, for each selected country, the observed value of the number of travelers to Italy and the one-step-ahead predicted levels in the out-of-sample period September 2012 - May 2019. The model seems to capture well the fluctuations of the phenomenon and the main turning points.
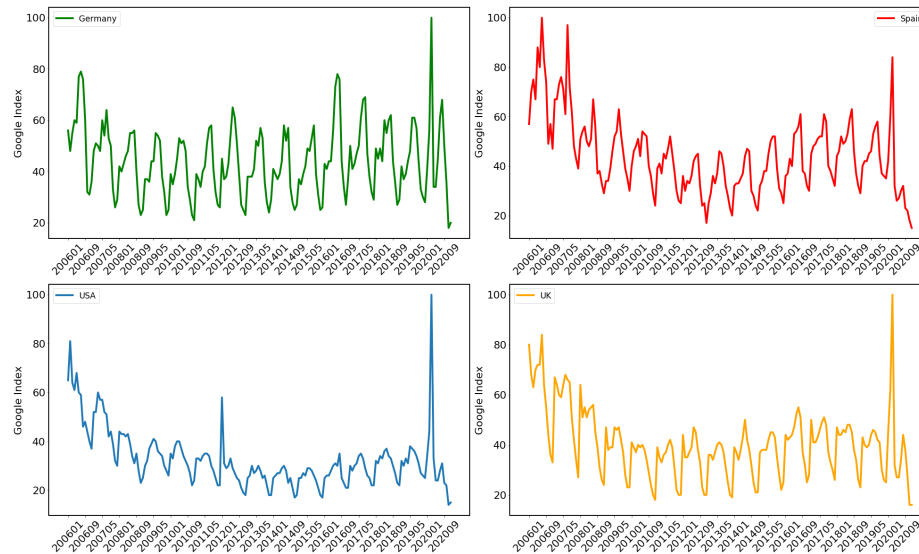


**Fig. 6.** Observed (solid line) and predicted (dashed line) number of travelers to Italy from Germany, Spain, United Kingdom and United States.

---

[25] The number of travelers at lag 1 and 12 is significant at 95% for all the considered countries.

[26] For Spain, the $R^2$ adjusted is 0.67 in the model without $GT$, and 0.77 in the one with the variable included. For the other countries, the $R^2$ is near to 0.9 in the model with only the AR component and the addition of the $GT$ index only increases it of around 0.01.

Although this data source proved to be interesting during the analyzed period, the Covid-19 pandemic pointed out its limits. Indeed, in March 2020 we witness a peak of search queries including the word "Italy" (see Figure 7 for selected countries), while in that month Italy was blocking the tourist inflow because of the pandemic. Such peaks very likely reflect the interest by the users in understanding the developing of the pandemic or in checking whether traveling to Italy was still doable or safe, even if they were not necessarily followed by actual travels.[27]

Therefore, in presence of extraordinary events, the Google classification seems to be less effective and the risk of outliers, given by false positive searches not related to tourism, increases significantly. Since the use of Google Trends strictly depends on the keywords included in the analysis, the use of other words as search queries, for example referring to specific Italian locations, could generate more accurate results.



**Fig. 7.** Search queries on Google including the word "Italy" by German, Spanish, UK and USA users.

## 6 Concluding remarks

In recent years, the Bank of Italy has carried out several experimental analyses to explore the possibility of integrating big data in the production of official

---

[27] This might explain why the peak in search queries appears also by considering only the *GT* "Travel" category .

statistics, in particular for compiling the "travel" item in the Balance of Payment.

The data that have been tested are appealing for their extraordinary timeliness and the amount of information offered, although they are very far away from being ready to use. Indeed, they need adjustments in order to define metrics that are coherent with the standards and the official definitions. The experiments often required adopting a trial-and-error approach to align these metrics to the prefixed standards and making strong assumptions that could potentially affect the results.

According to our tests, mobile phone data seem to be the most suitable ones to be integrated with the frontier survey, as they can produce reliable estimates of the number of international travelers crossing the Italian borders, thus potentially replacing, at least partially, the counting procedures in the Bank of Italy frontier survey. The Bank of Italy is already moving in this direction.

The other big data sources analyzed, electronic payments data and Google Trends data, showed more limitations and drawbacks.

Electronic payment data proved useful for achieving a preliminary estimate of total travel expenditures, as they are more timely than survey data. However, at this stage they can be used, at least from a BoP perspective, only for checking purposes. On the other hand, due to the informative potential of this source, we will continue exploring how to overcome the main problems by identifying the features that the data should have to be fully usable.

The Google Trends index proved to be useful for estimating the number of international travelers. But the sensitivity of the index to extraordinary circumstances like the Covid-19 pandemic needs to be further investigated before considering the integration of such an index in the compilation process.

## References

1. R. Ahas, A. Aasa, S. Silm and M. Tiru, 'Mobile positioning data in tourism studies and monitoring: case study in Tartu', Estonia,*Information and communication technologies in tourism 2007*, 2007, pp. 119-128.
2. R. Ahas, A. Aasa, A. Roose, Ü. Mark and S. Silm, 'Evaluating passive mobile positioning data for tourism surveys: An Estonian case study', *Tourism Management*, 29(3), 2008, pp. 469-486.
3. R. Ahas, J. Armoogum, S. Esko, M. Ilves, E. Karus, J.L. Madre, O. Nurmi, F. Potier, D. Schmücker, U. Sonntag and M. Tiru M, *Feasibility study on the use of mobile positioning data for tourism statistics*, Consolidated report Eurostat contract No 3051.2012.001-202.452, 2014.
4. G. Ardizzi, A. Gambini, A. Nobili, E. Pimpini and G. Rocco, 'L'impatto della pandemia sull'uso degli strumenti di pagamento in Italia', Banca d'Italia, Approfondimenti (Research papers), 8, 2021.
5. C. Artola, F. Pinto and P. de Pedraza García, 'Can internet searches forecast tourism inflows?', *International Journal of Manpower*, 36(1), 2015, pp. 103-116.
6. N. Askitas and K. F. Zimmermann, *Google econometrics and unemployment forecasting*, Technical Report, SSRN 899, 2009.

7. P. F. Bangwayo-Skeet and R. W. Skeete, 'Can Google data improve the forecasting performance of tourist arrivals? Mixed-data sampling approach', *Tourism Management*, 46, 2015, pp. 454-464.

8. Y. Carrière-Swallow and F. Labbé, 'Nowcasting with Google Trends in an emerging market', *Journal of Forecasting*, 32(4), 2013, pp. 289-298.

9. H. Choi and H. Varian, 'Predicting the present with Google Trends', *Economic Record*, 88, 2012, pp. 2-9.

10. J. M. Coelho, C. I. Ferreira and J. Veiga, 'The use of payments card data for travel statistics', Supplement to the statistical bulletin, Banco de Portugal, 2011.

11. F. D'Amuri and J. Marcucci, 'The predictive power of Google searches in forecasting US unemployment', *International Journal of Forecasting*, 33(4), 2017, pp. 801-816.

12. C. Demunter, 'Tourism Statistics: Early Adopters Of Big Data?', Sixth UNWTO International Conference on Tourism Statistics, 2017.

13. M. Dubreuil, 'The Potential Use of Credit/Debit Card Data for Tourism Statistics', UNWTO/DG GROW Workshop Measuring the economic impact of tourism in Europe: The Tourism Satellite Account (TSA), 2017.

14. P. Deville, C. Linard, S. Martin, M. Gilbert, F. R. Stevens, A. E., Gaughan, V. D. Blondel V. D. and A. J. Tatem, 'Dynamic population mapping using mobile phone data', *Proceedings of the National Academy of Sciences*, 111(45), 2014, pp. 15888-15893.

15. J. C. Deville and C. E. Särndal, 'Calibration estimators in survey sampling', *Journal of the American statistical Association*, 87(418), 1992, 376-382.

16. 'Decisions taken by the Governing Council of the ECB (in addition to decisions setting interest rates)', ECB's website, 11 December 2020.

17. 'Methodology for the compilation of international travel statistics', Estonia Central Bank's website.

18. 'Methodology and sources for compilation of the accounts of the Balance of Payments', Estonia Central Bank's website.

19. Google trends website, https://trends.google.com/trends/?geo=IT

20. IMF, *Balance of Payments and International Investment Position Manual*, Sixth Edition (BPM6), 2009.

21. A. Janecek, D. Valerio, K. A. Hummel, F. Ricciato and H. Hlavacs, 'The cellular network as a sensor: From mobile phone data to real-time road traffic monitoring', *IEEE transactions on intelligent transportation systems*, 16(5), 2015, pp. 2551-2572.

22. F. Le Gallo and K. Schmitt, 'Measuring international travel during the Covid-19 pandemic', Banque de France's website, post n 184.

23. J. Li, L. Xu, L. Tang, S. Wang and L. Li, 'Big data in tourism research: A literature review', *Tourism Management*, 68, 2018, pp. 301-323.

24. S. Lokanathan, G. E. Kreindler, N. N. de Silva, Y. Miyauchi, D. Dhananjaya and R. Samarajiva, 'The potential of mobile network big data as a tool in Colombo's transportation and urban planning', *Information Technologies and International Development*, 12(2), 2016, pp. 63-73.

25. N. McLaren and R. Shanbhogue, 'Using internet search data as economic indicators', Bank of England Quarterly Bulletin, 2011, Q2.

26. F. Ricciato, P. Widhalm, M. Craglia and F. Pantisano, 'Estimating population density distribution from network-based mobile phone data', Publications Office of the European Union, 2015.

27. T. Suhoy, 'Query indices and a 2008 downturn: Israeli data', Bank of Israel, 2009.

28. UN, *Manual on Statistics of International Trade in Services*, 2010.

29. UN Department of Economic and Social Affairs Statistics' website, 45th Session, chapter I, section B, 2014, p. 20.
30. UNWTO, 'International Recommendations for Tourism Statistics' (IRTS), 2008.
31. S. Vosen and T. Schmidt, 'Forecasting private consumption: survey-based indicators vs. Google trends', *Journal of Forecasting*, 30(6), 2011, pp. 565-578.
32. L. Yezekyan, 'Compilation of e-commerce data for Balance of Payments', IFC-CBA workshop on external sector statistics, 2018.