

Schmücker, Dirk; Reif, Julian

Article

Measuring tourism with big data? Empirical insights from comparing passive GPS data and passive mobile data

Annals of Tourism Research Empirical Insights

Provided in Cooperation with:

Elsevier

Suggested Citation: Schmücker, Dirk; Reif, Julian (2022) : Measuring tourism with big data? Empirical insights from comparing passive GPS data and passive mobile data, Annals of Tourism Research Empirical Insights, ISSN 2666-9579, Elsevier, Amsterdam, Vol. 3, Iss. 2, pp. 1-12, <https://doi.org/10.1016/j.annale.2022.100061>

This Version is available at:

<https://hdl.handle.net/10419/337720>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

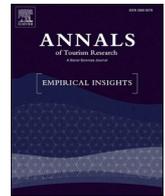
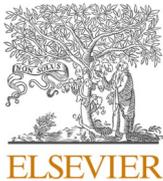
Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



<https://creativecommons.org/licenses/by-nc-nd/4.0/>



Measuring tourism with big data? Empirical insights from comparing passive GPS data and passive mobile data

Dirk Schmücker^{a,b,*}, Julian Reif^{a,2}

^a German Institute for Tourism Research at the West Coast University of Applied Sciences, Heide, Germany

^b Institute for Tourism Research in Northern Europe (NIT), Kiel, Germany

ARTICLE INFO

Editor: Dr. Kirilova Ksenia

Keywords:

Big Data
Mobile Network Data
Passive GPS Data
Spatio-temporal behaviour
Tourist classification

ABSTRACT

In this paper we aim to classify digital data sources for the measurement of tourist mobility, to establish a set of assessment indicators, and to compare two Big Data sources to gain empirical insights into how we can measure tourism with Big Data. For three holiday destinations in Germany, passive mobile data and passive global positioning systems (GPS) data are compared with reference data from the destinations for twelve weeks in the summer of 2019. Results show that mobile network data are on a plausible level compared to the local reference data and are able to predict the temporal pattern to a very high degree. GPS app-based data also perform well, but are less plausible and precise than mobile network data.

1. Introduction

The use of Big Data for measuring tourism has been discussed for some time (Ahas et al., 2014b). The actual number of tourists in a destination is crucial information against the backdrop of sustainable destination development. As such, great hopes are placed in Big Data to eliminate the shortcomings presented by official statistics (Demunter, 2017). Knowledge of the number of tourists is important, for example, since it has been shown that the life satisfaction of the residents of a tourist area is dependent on tourism intensity, among other factors (Tokarchuk, Gabriele, and Maurer, 2021). In addition, previous research has also established that a tourist's positive experience of a destination varies systematically with the number of tourists present (Tokarchuk, Barr, and Cozzio, 2021).

The term "Big Data" is of growing importance to tourism research (Mariani and Baggio, 2021). We identified three notions of what Big Data can mean in the context of this paper: First, the ubiquitous "Gartner's 3 V" of high-volume, high-velocity and high-variety (Laney, 2001) is a constituting element of Big Data (while other V's, added later, are not). Second, Big Data in the context of this study is data that has not been generated for the purpose of the study, but for some other reason. This aspect is somewhat pejoratively called "exhaust data" (Kitchin, 2014; O'Leary and Storey, 2022), although we argue that "upcycled

data" would be a better description. Finally, it seems that Big Data tends to be overestimated in terms of availability and validity. The perspective that Big Data is so big that no further care has to be taken for validity or reliability (Yang, Pan, Evans, and Lv, 2015) has been criticised by scholars in the past (Ferreira and Vale, 2020; Mazanec, 2020), and we agree with this scepticism.

However, Big Data sources are not only used for purely volume-based estimations of tourists in a destination; they are also especially used for measuring inter- and intraregional tourist spatio-temporal behaviour (Hardy, 2020). The knowledge of tourist movement patterns is so significant because it has an impact on the destination's infrastructural questions, transport system, product development, planning of new tourist attractions, and management of tourism-induced economic, ecological, and socio-cultural effects (Lau and McKercher, 2006).

We are currently witnessing a twofold digital transformation. First, the recent COVID-19 pandemic serves as a catalyst for the expansion of digital infrastructure in tourist destinations to manage tourism flows with sensors. Second, this investment in digital infrastructure has the potential, in turn, to transform a destination as a whole into a smart destination (Gretzel, 2018; Shafiee, Rajabzadeh Ghatari, Hasanzadeh, and Jahanyan, 2021).

Despite a wide variety of research on this topic, many questions

* Corresponding author at: German Institute for Tourism Research at the West Coast University of Applied Sciences, Heide, Germany.

E-mail addresses: dirk.schmuecker@nit-kiel.de (D. Schmücker), reif@fh-westkueste.de (J. Reif).

¹ NIT – Institute for Tourism Research in Northern Europe, Fleethörn 23, D-24103 Kiel, Germany.

² FH Westküste, Fritz-Thiedemann-Ring 20, D-25746 Heide, Germany.

remain unanswered. This can be attributed to rapid technological advances and the need to apply digital techniques to investigate further. The most pressing questions are those of validity and applicability: What do we actually measure with Big Data, and what does Big Data tell us about the real world? The validity of Big Data for tourism purposes, together with questions of applicability and accessibility, has been identified as a problem (Reif and Schmücker, 2020). As Big Data research in tourism usually examines only one data source with only partial insights (Park, 2021), combining methods and Big Data sources appears to be a possible solution to overcome those drawbacks and to validate Big Data with other data sources.

Against this backdrop, our main objectives in this paper are to assess how far selected Big Data sources can correctly identify main tourism segments (overnight visits and day trips) both in terms of volume and pattern over time. To achieve this objective, two data sources are contextualised within a conceptual framework of digital data sources for tourism measurement. Additionally, a set of quality indicators is introduced to allow for a systematic assessment. In distinction from other studies that rely solely on Big Data sources to describe real world-situations (Ramos, Yamaka, Alorda, and Sriboonchitta, 2021), this paper is innovative in two ways. First, we compare two Big Data sources, passive mobile data and app based data. Second, we use a destination perspective with reference data from the real world to validate digital data sources.

In this study we mainly compare two data sources suitable for measuring visitor frequencies and flows: passive mobile data and passive GPS data. Both sources rely on data generated by digital devices carried by the consumer. Additionally, two more data sources are used for reference: local reference data (from the destination management organization) and web-scraped data of bookings (from the Airbnb and the HomeAway platforms). For three holiday destinations in Germany, passive mobile data and passive GPS data are compared with reference data from the destinations for twelve weeks in the summer of 2019.

2. Literature review and conceptual framework

2.1. Using Big Data for detecting tourist spatio-temporal behaviour

The rapid advancements in technology and the need to understand the mobility of tourists both to the destination (inter-regional) and in the destination (intra-regional) within a framework of a sustainable destination development increases the possibilities in tracking people's spatio-temporal behaviour. For a better understanding of the spatio-temporal behaviour of tourists via technology, two research strands can be identified (Reif 2021):

(1) The use of experimental setups with sensors to capture the unaltered experience in the destination, usually with small sample sizes and with interdisciplinary research approaches (Bastiaansen, Oosterholt, Mitas, Han, and Lub, 2020; Birenboim, Dijst, Scheepers, Poelman, and Helbich, 2019; Reif and Schmücker, 2021; Scuttari, 2021; Shoval, Schvimer, and Tamir, 2018)

(2) The use of Big Data sources such as passive mobile data (Raun, Ahas, and Tiru, 2016), Twitter (Aagesen et al., 2020), or other forms of Big Data (Li, Xu, Tang, Wang, and Li, 2018)

Ignoring the promising approaches to understanding the emotional experiences of tourists with sensors (Sensing Tourists, Shoval, 2018), a multitude of digital data sources, most of them Big Data, are used to digitally track people in time and space or forecast tourism flows (Li, Zheng, and Ge, 2022; Yang, Fan, Jiang, and Liu, 2022). The approaches vary from passive mobile data (Saluveer et al., 2020; Zheng, Li, Lin, and Zhang, 2022) to Bluetooth (Versichele, Neutens, Delafontaine, and van de Weghe, 2012), Photo-Sharing Platforms like Flickr (Önder, Koerbitz, and Hubmann-Haidvogel, 2016), public Wi-Fi-networks (Ramos et al., 2021), destination cards (Zoltan and McKercher, 2015), Google Trends data (Höpken, Eberle, Fuchs, and Lexhagen, 2021), and bank card transactions (Aparicio, Hernández Martín-Caro, García-Palomares, and

Gutiérrez, 2021; Romero Palop, Murillo Arias, Bodas-Sagi, and Valero Lapaz, 2019), among others. However, the validity of Big Data can be a problem, as some instruments do not actually measure people or tourists but rather electronic devices, vehicles, or other non-personal or non-touristic objects (Reif and Schmücker, 2020; Schmücker and Reif, 2021). This ambiguity brings about the risk of being trapped into neo-positivist positions without theoretical reference by looking only at the data, to the exclusion of other potential influencing factors (Bauder, 2019).

Combining different Big Data sources seems to be a fruitful approach to detect both the strengths and weaknesses of the data, to understand possible biases (Nyns and Schmitz, 2022), to discover new insights, and to build tourism-oriented theories (Mazanec, 2020; Park, 2021). However, we can identify only a few papers in tourism research in which different Big Data sources are combined to estimate tourist spatio-temporal behaviour. For instance, Batista e Silva et al. (2018) combine data from Booking.com and TripAdvisor.com with data from official accommodation statistics to obtain the monthly tourist density grids in the EU. In another study, data from Panoramio, Foursquare, and Twitter are combined to reveal different tourist activities and hot spots in Madrid (Salas-Olmedo, Moya-Gómez, García-Palomares, and Gutiérrez, 2018). Nyns and Schmitz (2022) combine passive mobile data with scraped accommodation data from AirDNA and official statistics to estimate unobserved tourist overnight stays in Wallonia. Web-scraped Airbnb data was also used with official statistics from the United Nations World Tourism Organization (UNWTO) and World Travel & Tourism Council to map global tourist hotspots (Adamiak and Szyda, 2021).

2.2. Classification of tourist tracking techniques

Several categorisations of measurement methods for quantifying place use and tourist tracking have been discussed in the literature. A bibliometric meta-analysis reveals 31 different tracking techniques (digital and non-digital) resulting in six different categories for classifying tourist tracking techniques (Padrón-Ávila and Hernández-Martín, 2021). Apart from non-digital tracking techniques, Shoval and Ahas (2016) identified two different phases of digital tracking in tourism: (1) high-resolution GPS tracking (a mostly active method of handing out GPS sensors to the participant) vs. lower spatio-temporal resolution sources like passive mobile data, Twitter, etc. (a mostly passive method), and (2) a smartphone revolution that makes use of all the sensors built into high-end smartphones. Hardy (2020) gives an overview of different state-of-the-art tourist tracking techniques with a pronounced focus on the ethical questions of each of the seven discussed tracking solutions. From a geographic perspective, Reif (2021) proposes a classification of digital data sources based on their spatial coverage on the macro, meso, and micro level.

2.3. Conceptual framework

In the conceptual framework presented here, we identified four categories of measurement methods for the quantification of place use (How many people are on site right now?), the quantification of mobility or tourism flows (What are the visitor flows in a certain area?), the identification of activity spaces (What movement range do people have?), and the origin-destination relations (Where do people come from and where do they go?), each with several sub-categories (Fig. 1).

Multi-spot measurements (Category A) review large areas and are usually filtered to identify signals in the area of interest. These measurements are not restricted to small locations (i.e., spots), but rather they extract location data from larger datasets that frequently contain signal or mobility chains. These chains make it possible to infer relevant information from these sources, including the origin and destination or the activity spaces of signals. This data can be obtained from devices or from the network infrastructure. GPS-based location data is obtainable

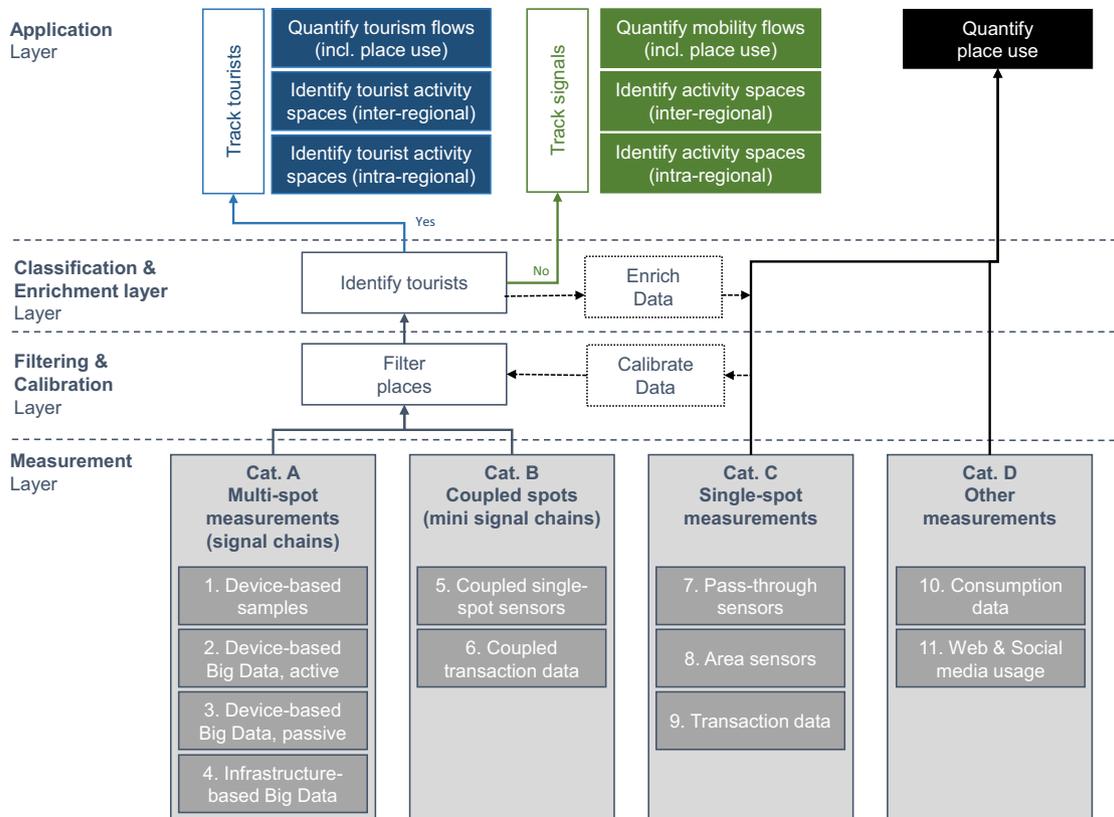


Fig. 1. Measurement categories and layers in measuring mobility and tourism. Source: Authors.

from mobile apps (Möhring, Keller, Schmidt, and Dacko, 2021) or car navigation systems and are stored in the device. Conversely, data from cell phone towers are stored in the network, not in the device (Raun et al., 2016). These are examples of Big Data, as opposed to sampled data, which can be collected by sampling visitors (e.g., active GPS tracking) (Bauder and Freytag, 2015). When multi-spot measurement data is collected over a long period, it is possible to make an educated guess about the usual environment of a signal, which is one basis for categorising tourists and non-tourists based on international conventions, such as from the United Nations Department of Economic and Social Affairs and Statistics Division (UNSD) and the United Nations World Tourism Organization (UNWTO, 2010). This categorisation occurs in the classification layer.

In the application layer, in theory, multi-spot sensors can be used to identify tourists and then to quantify tourism flows and tourism action spaces with intra- and interregional movement patterns. Without the classification layer, multi-spot sensors can be used to quantify mobility flows, activity spaces, and origin-destination relations (not tourist-specific). In addition, quantifying place use is possible by filtering the multi-spot dataset to the spot or area of interest.

In Category B, we find coupled transaction data and coupled sensors. Transaction data (e.g., from a destination guest card system) (Zoltan and McKercher, 2015) typically consists of one-spot measurements, but card holders can be re-identifiable at all spots where they present their personalized guest card. The same principle can be found in coupled data from Wi-Fi (Yamamoto, Sato, and Kamitani, 2021) or Bluetooth (Versichele et al., 2012) sensors. Typically, these trackers store the MAC address of the devices they track, and, through a combination of various sensors and re-identification of MAC addresses, they can establish local mobility chains. Compared to the types of Big Data in Category A, the data in Category B can be described as “small Big Data” because the resulting signal chains are considerably shorter and restricted to selected

points.

Category C contains one-spot measurements, which can be provided via pass-through sensors (e.g., photo sensors, laser scanners, induction loops) or area sensors (e.g., LiDAR sensors, smart cameras, or CCTV) (Izquierdo Valverde, Prado Mascuñano, and Velasco Gimeno, 2016). One-spot (i.e., “single-spot”, Hardy, 2020, p. 62) measurements are restricted to one spot for each sensor, that is, small localities or points of interest such as entrance areas, parking spaces, or spots on hiking or biking trails. In practice, a point of interest can have one measurement spot (e.g., the entrance of an outdoor information centre) or more than one measurement spot (e.g., several access roads to a parking area), which then need to be matched to obtain correct volume numbers. Usually, there are only limited possibilities to infer additional information from the data. Some sensors (e.g., some types of laser scanners) can measure people’s height and thereby classify children and adults. Other sensors (e.g., two or more pyroelectric sensors installed at different altitudes along a track) can measure the height of an inferred source to classify walkers, bikers, and horse-riders. One-spot sensors are usually very precise, but they are also restricted in what information they can provide about the signal. Applications of single-spot measurements therefore mainly serve to quantify place use. A third subcategory of Category C is transaction data, which is produced whenever a transaction occurs (e.g., the purchase of entrance or parking tickets or the use of a cash-dispenser).

Category D contains other forms of measurements, including consumption data (e.g., inferring information on place use from the amount of electricity or water used in the area) and web and social media data (e.g., inferring information on place use from the pictures of the place shared on a social media platform, such as ambient geospatial information). The latter method plays a particularly important role in tourism research with data obtained from Twitter (Hawelka et al., 2014), Flickr (Girardin, Calabrese, Fiore, Ratti, and Blat, 2008), Weibo

(Chen, Becken, and Stantic, 2021), or Instagram (Paül i Agustí, 2020). Theoretically, this data source also can be used to identify local mobility chains, which would possibly place it in Category B. However, data from this source is usually very limited in terms of completeness and validity, as has been shown with photos posted on social media (Bauder, 2019).

The measurements in Categories A and B produce signal chains, which can in turn be used to quantify place use and to identify inter-regional activity spaces and intra-regional activity spaces. Measurements in Categories C and D can usually only quantify place use but with relatively high precision. Therefore, data from Categories C and D can also be used to calibrate data from Categories A and B. In the opposite direction, data from Categories C and D can be used to enrich data from Categories A and B.

For measurements in Categories A and B (but also for some in Category C, such as cameras and Wifi/Bluetooth sensors), ethical issues (Hardy, 2020) and data-protection issues can—and in many cases will—occur. In Europe, pictures of people, car registrations, and identifiable device properties, like the device-ID, Ad-ID, MAC-address, or Digital Subscriber Number, are frequently considered “personal data” and fall under the General Data Protection Regulation (Regulation (EU) 2016/679) (Ahas et al., 2014a; Reif and Schmücker, 2020). The use of data from these sources is therefore restricted (e.g., there are requirements of pseudonymisation or anonymisation and/or the explicit consent of users to use their data).

2.4. Assessment dimensions

To assess the qualities and applicability of the data sources, we suggest a set of 13 indicators (Table 1). The first three are tourism-specific and are as follows:

- (1) The external validity (i.e., are we measuring people or something else?)
- (2) The ability to derive object properties in addition to simply counting signals (e.g., movement directions, movement paths, bike vs. walk, tall persons vs. small persons)
- (3) The ability to classify important tourism segments (e.g., overnight visits, day trips, trips of local inhabitants, and other types of mobility)

Other (non-tourism) types of mobility are manifold and include job and school commuting; trips for private medical, sports, or administrative reasons; and professional mobility like the movements of taxi, bus, or lorry drivers, deliveries, and police, fire or ambulance services.

The next group of dimensions is related to space and time. Spatial granularity describes the resolution of signals in three-dimensional space. Temporal granularity refers to the resolution of signals over time, including the question of whether signals are produced in fixed intervals (e.g., a camera picture every 5 min) or are triggered by events (e.g., the passing of a light sensor). Latency is the second time-related dimension, and it refers to the time before the data can be used (lead

Table 1
13 assessment dimensions.

Specific touristic dimensions	(1) External validity (2) Object property detection (3) Tourist classification
Time & space	(4) Spatial granularity (5) Temporal granularity (6) Latency (time until use)
Generic dimensions	(7) Completeness (8) Precision of signal identification (9) Reliability (number of dropouts)
Social/organisational dimensions	(10) Accessibility and cost (11) Transparency of data processing (12) Compliance with data protection rules (13) Ethical justifiability

time). Latency, therefore, is the time between measurement and availability of data for the researcher.

More generic data quality indicators include completeness, precision, and reliability. Completeness refers to the amount of signals a sensor measures in relation to all the signals produced. An example of completeness is how mobile network data usually provides the mobile signals from only one network, rather than all mobile signals. Precision refers to the share of signals in relation to all the signals to be measured by the sensor. For example, laser scanners have a 98% precision, meaning that they miss 2% of all the signals they should measure. Reliability is measured against the fact that all technical systems can have signal drop-outs because of downtimes, technical issues in transferring data, and so on.

The last category addresses social and organisational dimensions. First, there is the accessibility and the cost of data. The production of sensor data is costly (such as in the case of infrastructure-based data) and/or will not be shared for privacy reasons (such as in the case of transaction data). Transparency in data processing can also be an issue if the organization does not control the data stream from the sensor to the data hub and must buy data from third parties, as is frequently the case for Big Data sources in Category A. Finally, compliance with data protection rules and ethical justifiability are relevant assessment dimensions. We will use these indicators in Section 5.

3. Methodology

3.1. Study area

Data was analysed for three destinations on the German North Sea coast and one additional attraction. These include Biusum and St. Peter-Ording, which are municipalities on the mainland, and Amrum, an island (Fig. 2). Table 2 shows data for the three places, as available from official accommodation statistics. All three destinations are located in the Wadden area of the North Sea; thus the marine landscape and a variety of water-related activities are the main tourist attractions. In addition, the Wadden Sea is a United Nations Education and Cultural Organization (UNESCO) World Heritage Site. The tourism intensity is far above the German average of 6.4 nights per inhabitant. All three destinations are heavily dependent on the domestic market. According to official accommodation statistics, about 98% of nights sold in the area are from the German source market.

We also analysed data for the Multimar Wattforum, an indoor nature information and exhibition centre in the same region. Overnight tourists do not play a role in this destination, as the centre does not provide any accommodations. Because the spot has a very limited spatial extension, we could only use GPS-based data for analysis. Mobile network data would have been far too coarse for this spot.

The data used for this study spans 12 weeks in the summer of 2019, starting in calendar week 23 (3 June) and ending in calendar week 34 (25 August). We chose the year 2019 to use data free from COVID-19-related distortions. As the pandemic temporarily extinguished tourism activity in the three places, no meaningful analysis would have been possible. We chose the beginning of the summer season because, in the summer, we see high-volume demand. We also chose this time frame because the school holidays provide a natural source of variation in the data, as the start date of school holidays in Germany differs from state to state. Although we would have preferred to have a longer time series, more temporal granularity (e.g., per hour or day instead of per week), and a larger sample, we conceded that access to most of the data used in this study is still not only a methodological, but also a financial and organisational challenge (Schmücker and Reif, 2021). Details for each data source are discussed in the following sections.

3.2. Data sources

In this study, we use multi-spot measurements based on the devices



Fig. 2. Study area.
Source: Authors.

Table 2
Study area.

Place	Büsum	Amrum	St. Peter-Ording
Type	Municipality	Island	Municipality
Inhabitants	4907	2263	3997
Accommodations establishments	171	188	157
Beds	8677	6188	10,896
Campsites	3	2	?
Pitches	460	n/a	663
Arrivals	285,397	101,100	379,908
Bednights	1,431,231	858,808	1,853,254
Bednights per inhabitant	291.7	379.5	463.7

Source: Statistical office for Hamburg and Schleswig-Holstein, own calculations.

carried by the consumer. Location data from a mobile app is in the device-based category (i.e., location data is stored in the device), while cell tower data from mobile network operators falls into the infrastructure-based category (i.e., data is stored in the network).

Both the data sources studied are able to classify signal chains into

overnight stays, day trips, trips made by inhabitants, and other forms of mobility. The classification is performed for each day separately and always refers to the place under observation. Classifications are based upon the identification of home and work zones. For example, an overnight tourist in Büsum is a person who spent the night before or after the day under observation in Büsum and whose residence is not in Büsum. A commuter to St. Peter-Ording is a person that has his or her workplace in St. Peter-Ording, but the place of residence outside of St. Peter-Ording. A day trip to Amrum is a person that has neither the place of work nor the place of residence on the island and which did not spend the night before or after the day under observation on Amrum. The details of the classification algorithms are not known and are kept secret by the data providers. Therefore, we assess the quality of classification by examining the results.

3.2.1. Infrastructure-based Big Data: Passive mobile network data

Passive mobile network data, also called (passive) mobile positioning data (Ahas, Aasa, Roose, Mark, and Silm, 2008; Nilbe, Ahas, and Silm, 2014), passive mobile data (Reif and Schmücker, 2020), or mobile

phone tower tracking (Hardy, 2020), is data derived from mobile phone signals in the networks of mobile network operators (MNO). This data is labelled “passive” because it occurs without any activity on the part of the user (besides switching the device on and letting it connect to the network). Mobile devices connect to the network either in regular intervals or, much more important today, when an *event* occurs. Events include handling phone calls, text messages, or internet connections (the most frequent use today) as well as changing cells (Janevski, 2019, p. 137 ff.; for the situation in Germany, see Sauter, 2018). In contrast to device-based GPS data, this data is infrastructure-based and only exists in the network, not in the device. Passive mobile data can be used for tourism statistics (Demunter, 2017); however, at least from a German perspective, strict data protection rules used to make it impossible to distinguish between tourist and non-tourist signals (Reif and Schmücker, 2020). This restriction has only recently been overcome.

We used data from the German network of Telefónica A.S. In Germany, Telefónica’s mobile phone services are marketed under the main brand O2 and several secondary brands. In the second quarter of 2019, the company reported 43.2 million active SIM cards, which represents a market share of 31.6%, with Deutsche Telekom (32.7%) and Vodafone (35.7%) being its only competitors (Bundesnetzagentur, 2021). Data was provided by Teralytics AG, a firm located in Zurich, Switzerland, specialising in passive mobile data. Teralytics uses its own system of rectangular geogrids, and data is available for these grids only. Because cell towers do not physically cover rectangular areas, the grids are the result of internal signal recalculations. The data provider neither reveals the algorithms for these recalculations nor shows the exact mechanics of identifying different target groups and projections to population level. Contrasting with earlier data (Reif and Schmücker, 2020), the provider can now identify home locations and work locations and infer tourist and non-tourist activity from the data. For Büsum and Amrum, the Teralytics geogrids cover the whole municipality. However, in the case of St. Peter-Ording, the grid configuration does not correspond well with the municipality borders. Therefore, we used mobile network data only for Büsum and Amrum. To identify tourist signals, regularities in movement patterns based on the home and work locations are analysed by the MNO, albeit with all the pitfalls described in Reif and Schmücker (2020).

3.2.2. Passive device-based Big Data: Passive GPS data

Compared to GLONASS, BeiDou, and Galileo, GPS is the most used variant of global navigation satellite systems (Chen et al., 2021). It is based upon a set of satellites which send their position and a timestamp in regular intervals. Receivers use a subset of satellites (at least three or four, depending on data requirements) and compute their own location relative to that of the satellites. Precision is usually in the range of several meters or better (Schaefer and Pearson, 2021).

Passive GPS data was provided by Meteonomiqs, a provider that uses location events from the smartphone application *wetter.com* for business optimization and forecasting, particularly in the retail sector. However, data from the location events can also be used for tourism purposes, as we demonstrate here. According to Meteonomiqs, their application *wetter.com* has more than 2.5 million users (over 90% of whom are Android users) who produce more than 700 million geo-location events per month. A location event is generated when a user accesses the app in a designated area, that is, a destination that can be defined using official statistical boundaries or geofencing zones. The user must give their consent to use the GPS-signal in the smartphone while using the app.

To identify tourist signals, historical data is used to identify regularities in the movement patterns. Together with the provider, we define tourist signals based on the place of residence. The place of residence is determined based on the most frequent place of overnight stay, with a rolling window of 50 days. After 26 successive overnight stays, the overnight stay location is marked as the place of residence. Aggregation is based on the corresponding geoshapes and must have a minimum of five users according to the General Data Protection Regulation (see

Section 3.1). Within the geoshape, the signal is classified as “resident”. Outside of the geoshape, the signal is classified as “commuter” or “tourist”. To address the concept of the “usual environment”, according to the UNSD and UNWTO (2010), visits of ten or more weeks per quarter are classified as “non-tourist mobility” (i.e., commuters) and visits for less than ten weeks per quarter in the corresponding geoshape are classified as “tourists”.

3.2.3. Scraped platform data

We also utilise the overnight volume data scraped from the distribution platforms of Airbnb and HomeAway. Such data has been used before in tourism research and in applied tourism statistics (Agarwal, Koch, and McNab, 2019; Gibbs, Guttentag, Gretzel, Yao, and Morton, 2018; Leick, Kivedal, Eklund, and Vinogradov, 2021). The data was scraped, processed, and delivered by a commercial data provider, AirDNA LLC (Denver, USA). This data does not represent anything near the true overall tourism volume in the destination, as only two distribution channels are included. Clearly, this data can also be unrelated to tourist movements. Therefore, we used this data source as an additional reference for the temporal distribution pattern analysis only.

3.2.4. Local reference data

Local reference data for overnight tourism was obtained from the three local tourism administrations. In all three places, tourists must pay a local tourism tax for each night spent in the location. Accommodation establishments collect the tax and report the daily number of guests to the local tourism administration. This number is known to be lower than the true volume of tourism demand for two reasons. First, business guests and tourists in owned dwellings or in accommodation provided without charge by relatives or friends are exempt from this type of tax. Second, the accommodation establishments might not collect or report the tax, even though this constitutes fraud. However, the volumes reported here are considerably higher compared to official accommodation statistics, which only includes accommodation establishments and campsites of ten and more beds or pitches.

In the case of Büsum, we also used data from the tax tickets for day visitors. Legally, any day visitor is obliged to obtain a day tax ticket, but it is well known that only few do—consisting mostly of those who want to go to the beach and must pass a gate. Therefore, this data source also underestimates the true volume of day visits but can be used for temporal pattern analysis. For Multimar Wattforum, we used transaction data from the centre’s ticket counter. This data reflects the true number of visitors to the centre with only a very small margin of error (such as from free admissions). Statistical data on the number of inhabitants is only available for a whole quarter; thus this data has the same averaged value for each calendar week. We therefore used this data only to establish a maximum line of reference.

3.2.5. School holidays

We hypothesized that school holidays could have a considerable influence on the tourism demand in the three destinations. As shown above, all three destinations are heavily dependent on the domestic market. We therefore restricted our analysis to Germany. In Germany, school holidays differ from state to state. We calculated the share of the German population living in states with school holidays for every day and averaged the result over each of the twelve weeks. Fig. 3 shows that, after a first peak in week 24 (the week after Pentecost), the curve starts to build up in week 26 to reach a peak in week 31, when all German states had school holidays.

3.3. Data analysis

Data analysis was performed in two steps. For the first step, we conducted a graphical exploratory data analysis showing the measured values over the course of the twelve weeks in the temporal scope. In the second step, we performed a statistical analysis. In this step we used a

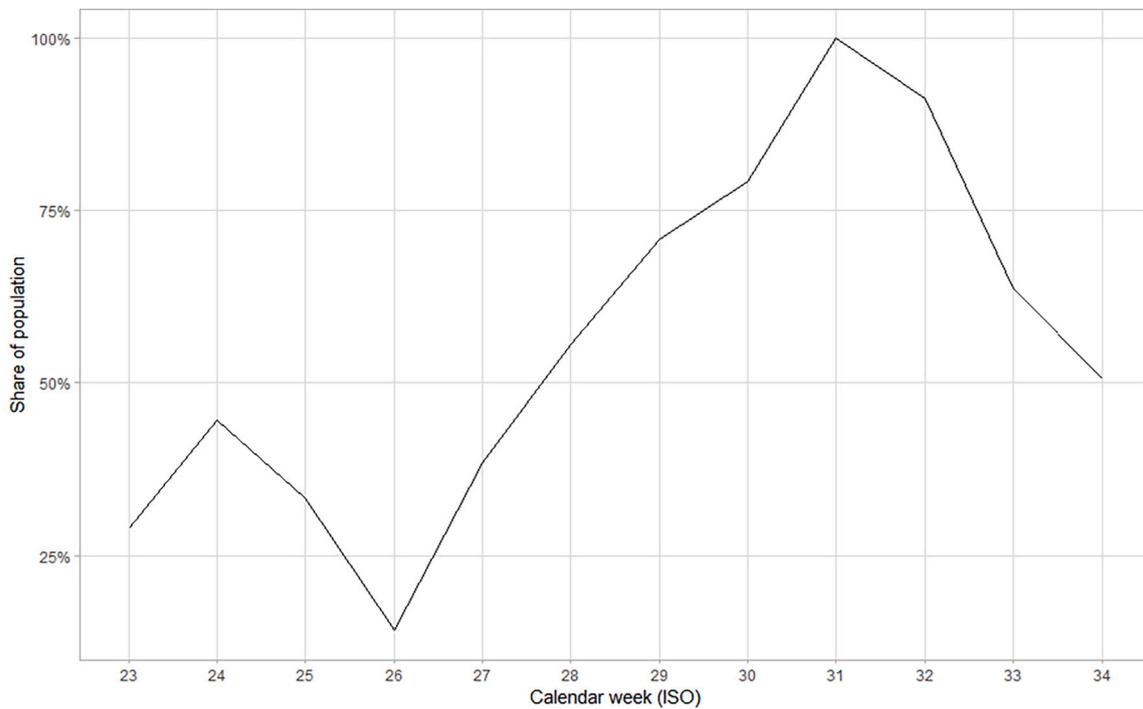


Fig. 3. Share of German population with school holidays (weekly sums, summer 2019). Source: Authors.

distance metric to assess the volume level and a correlation coefficient to assess the pattern over time in relation to the reference data. To measure the volume levels, we used the sum difference over the twelve weeks and the absolute distance of two vectors (L1 norm or “Manhattan” distance, Minkowski, 1910, p. 2), as implemented in the *dist* function of the R base package. Using only the sum difference could lead to misinterpretations, as the sum over twelve weeks might be identical for two vectors, although there is a considerable difference in every week (e.g., vector 1

is below vector 2 half of the time and above vector 2 the other half of the time).

A correlation coefficient was used to reflect the degree of the linear relationship between the two vectors, thus reflecting the similarity of the distribution pattern over time. We used Pearson’s *r* for this part of the analysis, as implemented in the R *psych* package (Revelle, 2021). However, the correlation coefficient is insensitive towards the level of the data. As long as the pattern shows that data points are above or

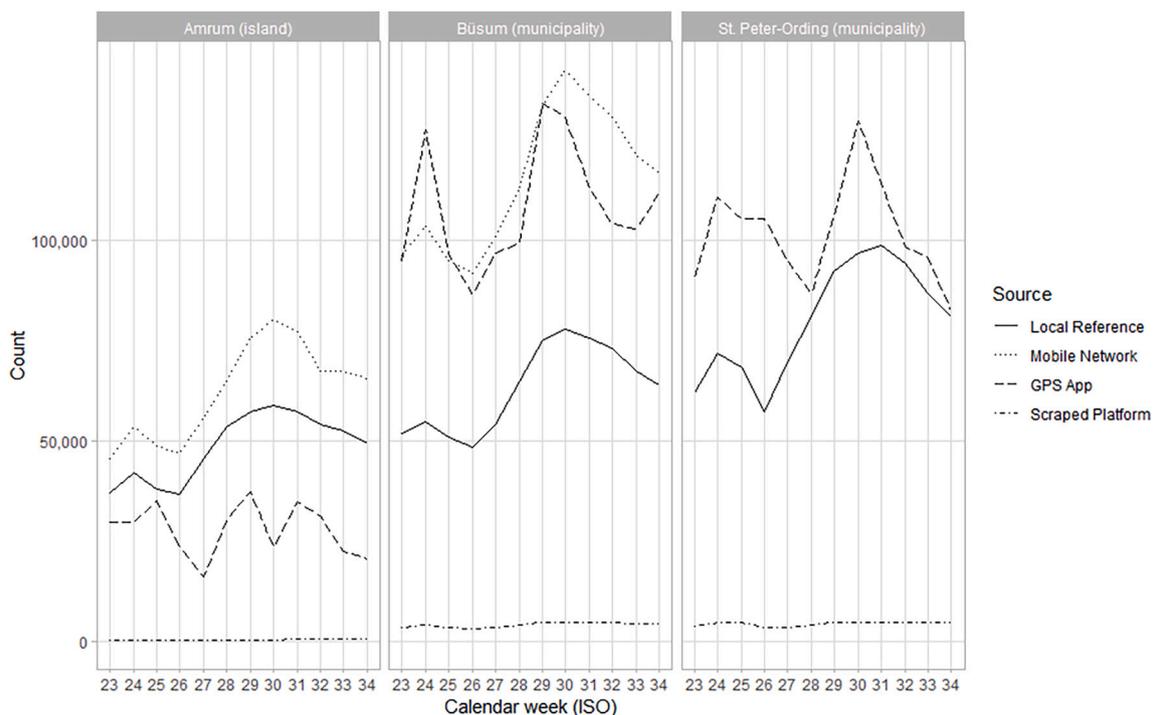


Fig. 4. Overnight tourists, by place and source (weekly sums, summer 2019).

below their respective average to the same degree, the coefficient will be high, even if the data points are on very different levels and thus very far apart. A correlation matrix cannot be used to infer predictions or to measure non-linear relationships (Ferber, 1956). However, the limited number of data points (12 weeks per place per segment per source) does not allow for the usage of more advanced machine learning algorithms, such as decision tree regressors (Witten, Frank, and Hall, 2011), as implemented in the Predictive Power Score (van der Laken, 2021; Wetschoreck, 2020) or other advanced algorithms. We did not calculate any inferential statistics (*p*-values or Bayes factors) because the setup of this study does not include any aspect of randomness or chance which would allow us to infer from the data to a population or between groups.

4. Results

4.1. Overnight tourists

Fig. 4 shows that for overnight tourists, the mobile network data shows a slightly, as is the case for Amrum, or considerably, as is the case for Büsum, higher volume level compared to the local reference data. For GPS app data, the situation is less clear. Although the volume reported by this data source is smaller in Amrum than the local reference data, the volume reported is greater in Büsum and in St. Peter-Ording. As expected, scraped platform data is on a considerably lower level than all other data sources.

As for the pattern over time, we see a first smaller peak in week 24, a steep rise in demand starting in calendar weeks 26 through 28, and a peak in weeks 30 through 32. The patterns are quite similar for all data sources, although GPS app data shows an unexpected plunge in Amrum in week 30 and in St. Peter-Ording in week 28.

To assess the differences in volume, we analysed the volume in relation to the local reference data and the L1 distance over the 12 weeks in relation to the local reference volume (Table 3). Results show that mobile network data is 29% (Amrum) and 82% (Büsum) higher than the local reference. The L1 index reveals approximately the same proportions between Amrum and Büsum, with Büsum having the higher deviation compared to Amrum. (See Tables 4 and 5.)

GPS app data is 71% percent higher in Büsum and 27% higher in St. Peter-Ording, but the data is 42% lower in Amrum (index 0.58). The L1 index shows that, again, Büsum has the highest deviation from the local reference data. As expected, scraped platform data has the lowest level and, therefore, the highest deviation from the reference data.

Analysing the patterns over time reveals very high correlations (author's own rating following Cohen, 1988) between the local reference data and the mobile network data and moderate or low correlations between the local reference data and the GPS app data. Correlations

Table 3
Overnight tourists, level of demand by data source and place.

Overnight tourists	Absolute volume over 12 weeks	Volume in relation to local reference (index)	L1 distance to local reference	L1 Distance in relation to local reference volume (index)
Local reference	A: 582,047 B: 758,003 S: 959,598	A: 1.00 B: 1.00 S: 1.00	A: 0.00 B: 0.00 S: 0.00	
Mobile network	A: 747,976 B: 1,380,189 S: n.a.	A: 1.29 B: 1.82 S: n.a.	A: 179,756 B: 674,035 S: n.a.	A: 0.31 B: 0.89 S: n.a.
GPS app	A: 334,485 B: 1,297,561 S: 1,220,555	A: 0.58 B: 1.71 S: 1.27	A: 268,192 B: 584,521 S: 282,703	A: 0.46 B: 0.77 S: 0.29
Scraped platform	A: 4098 B: 49,188 S: 52,198	A: 0.01 B: 0.07 S: 0.05	A: 626,111 B: 767,883 S: 983,017	A: 1.08 B: 1.01 S: 1.02

Note: A: Amrum (island), B: Büsum (municipality), S: St. Peter-Ording (municipality); n.a.: not available.

Table 4
Results (Pearson's r) for overnight tourists.

Overnight tourists	Local reference	Mobile network	GPS app	Scraped platform
Mobile network	A: 0.981 B: 0.994 S: n.a.			
GPS app	A: 0.123 B: 0.639 S: 0.347	A: 0.134 B: 0.681 S: n.a.		
Scraped platform	A: 0.662 B: 0.965 S: 0.781	A: 0.662 B: 0.962 S: n.a.	A: - 0.107 B: 0.723 S: 0.295	
Share of population with school holidays	A: 0.897 A: 0.939 S: 0.971	A:0.891 B: 0.934 S: n.a.	A: 0.324 B: 0.550 S: 0.344	A: 0.700 B: 0.933 S: 0.741

Note: A: Amrum (island), B: Büsum (municipality), S: St. Peter-Ording (municipality); n.a.: not available.

Table 5
Results for day trips.

Overnight tourists	Local reference	Mobile network	Absolute volume over 12 weeks	Volume in relation to local reference (index)
	Correlation		Volume	
Local reference			A: n.a. B: 49,714 M: 77,510	A: n.a. B: 1.00 S: 1.00
Mobile network	A: n.a. B: 0.934 M: n.a.		A: 93,957 B: 374,592 M: n.a.	A: n.a. B: 7.53 M: n.a.
GPS app	A: n.a. B: 0.891 M: 0.874	A: 0.721 B: 0.896 M: n.a.	A: 59,485 B: 555,444 M: 112,737	A: n.a. B: 11.72 M: 1.45

Note: A: Amrum (island), B: Büsum (municipality), M: Multimar (attraction); n.a.: not available.

between the local reference data and the scraped platform data are moderate or (very) high. There is also a high correlation between the reference data and the share of population having school holidays. Therefore, school holidays can be used to predict 80% (R^2) or more of the variation in tourism demand in a simple linear regression model (Cohen, 1968).

4.2. Day trips

Fig. 5 shows the distribution patterns for day trips. In Amrum, the numbers are much smaller than in Büsum, which might be explained by the fact that Amrum is an island and can only be reached by ferry. Mobile network data in Amrum is somewhat higher than the GPS app data. In Büsum, mobile network data is considerably lower than the GPS app data.

As expected, local reference data in Büsum is much lower than the two other sources. For the Multimar Wattforum attraction, local reference data can be considered very reliable, as it is derived from the transactions at the centre's ticket counter. GPS app data is somewhat higher than the reference. The two lines are very close in the first six weeks, but they diverge in the second half of the period analysed.

Correlations are high for all the combinations analysed but are lowest (0.874) for Multimar, where the local reference data is most reliable. In terms of volume, we can compare reference data and GPS app data for the Multimar attraction. The GPS app data is about 45%

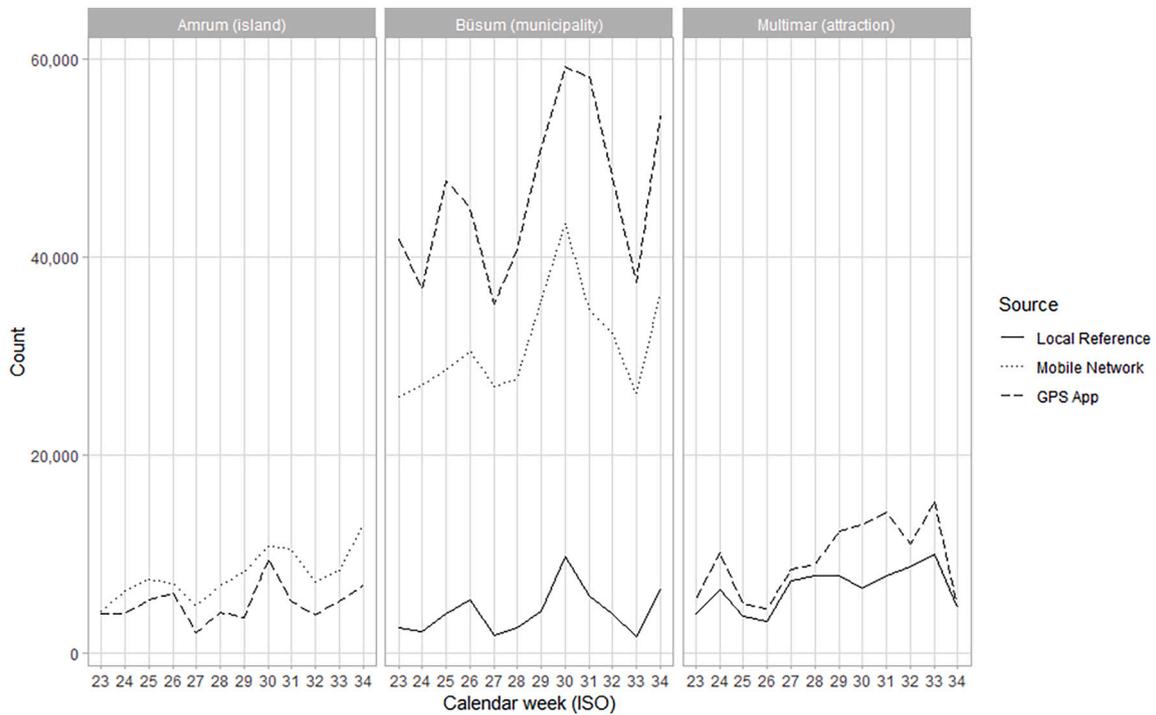


Fig. 5. Day trips, by place and source (weekly sums, summer 2019).

higher compared to reference data from the ticket counter.

4.3. Local population (inhabitants)

Local population data is available for Amrum, Büsum, and St. Peter-Ording. As outlined above, weekly reference data is not available; thus, the reference line in Fig. 6 is a straight line for all three destinations. It is plausible to expect that the number of inhabitants in the spot is usually smaller than the registered population, specifically during school

holidays (which start in week 27 and end in week 32 in the state of Schleswig-Holstein). Fig. 6 shows that the mobile network data and the GPS app data are below the reference line in most cases and never exceed the reference line. However, the GPS app data touches the reference line in eight out of the twelve weeks for Amrum, one out of the twelve weeks in Büsum, and two out of the twelve weeks for St. Peter-Ording.

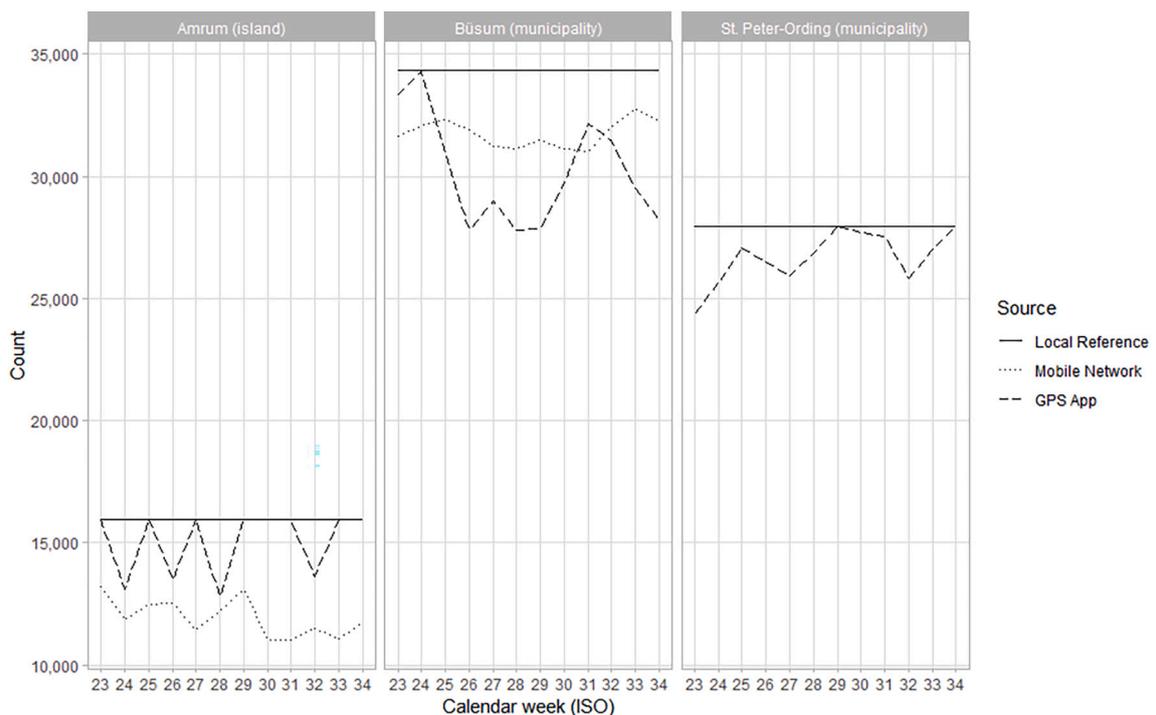


Fig. 6. Local population, by place and source.

5. Discussion

Results for overnight tourists are, in most cases, higher than the local reference data. This is plausible because the local reference is known to underestimate the true overnight volume. However, in the case of Amrum, the GPS app data is lower than the local reference over the whole period. This result is not plausible. It is not clear if this result is due to the spatial location (island), weather-effects, or systematic errors (issues with algorithms). Whether or not the mobile network data and the GPS app data actually show the true volume of overnight guests cannot be assessed based upon the reference data available for this study.

As for the distribution pattern over time, the mobile network data show a very high correlation with the local reference data. In both cases studied, the correlation is very close to 1, leading to a predictive value in a linear model (R²) of close to 100% (96% and 99%). However, the correlation coefficients for the GPS app data are considerably lower in all three cases studied. We cannot assess whether this is a systematic issue with this category of data (Category 3: Device-based passive Big Data) or an issue with this individual data source. Further assessment and replication studies with data from different providers are needed.

When comparing both data sources, it might be worth noting that the destinations in this study have an overwhelming proportion of national guests. Therefore, mobile network data in our study can bypass some complexities of international visitors using local networks. These complexities result from the fact that MNO can only record what is happening in their own network. The visibility of international visitors in a given network depends on the roaming contracts between the home MNO and the MNO in the visited country. In addition, residents of one country may use MNO contracts from another country. As a consequence, the international MNO user base can be volatile, thus leading to noisy data.

For day trips, the analysis is more incomplete. Local reference data and data from the mobile network and the GPS app were available for only one place. In terms of correlation over time, the mobile network data reaches a somewhat higher correlation (0.934) with the local reference data than with the GPS app data (0.891). However, the difference is moderate. In the case of the Multimar attraction, only the GPS app data was available. This data showed a high correlation of 0.874 with the local reference data, but the volume was overestimated by about 45%.

For inhabitants, data are plausibly below the reference line. However, in the case of GPS app data, there are several weeks when the data just touches the reference line. This might be an indicator that data have been manually cut off by the data provider in order to avoid exceeding the reference line. However, due the intransparency of the data preparation process this is only an assumption.

One objective of this study was to further assess the quality and the applicability of the selected data sources for measuring and tracking tourism activities along the lines of the assessment dimensions introduced above. Table 6 outlines the complete assessment based on the conceptual framework and the empirical results of this study (for an explanation of each assessment dimension please refer to Section 2.4 of this paper).

6. Conclusion

It has become clear that an assessment of volume is particularly challenging due to a lack of reliable reference data. This is true for overnight visits and even more true for day trips. However, it would be premature to simply declare the mobile network and the GPS app data as the new reference. One way out of this dilemma could be to use spot measurements (Category C) to calibrate the multi-spot measurements used here.

We recommend using multi-spot measurements (Category A), specifically mobile network data, to assess the distribution patterns over

Table 6

Assessment results of passive mobile data and passive GPS data.

Specific dimensions	Passive mobile data	Passive GPS data
(1) External validity	Limited (measures devices, not persons)	Limited (measures devices, not persons)
(2) Object property detection	Numerous (large signal chains)	Numerous (large signal chains)
(3) Tourist classification	Reasonable	Less reasonable
<i>Time & space</i>		
(4) Spatial granularity	Limited (restricted to cell towers and analytical cells)	High (several meters or better)
(5) Temporal granularity	Adequate, but frequency of events varies	Adequate, but frequency of events varies
(6) Latency (time to use)	Days or longer	Dependent on data provider
<i>Generic dimensions</i>		
(7) Completeness	Limited, as one network usually covers only a portion of the population; user base might be lopsided	Limited, but highly dependent on data provider; user base might be lopsided
(8) Precision of signal identification	High	High
(9) Reliability (number of drop-outs)	Moderate, occasional drop-outs occur	Unclear
<i>Social/organisational dimensions</i>		
(10) Accessibility and cost	Quite costly, small number of providers	Moderate, growing number of providers
(11) Transparency of data processing	Low	Low
(12) Compliance with data protection rules	High, but needs expensive anonymisation procedures	High (app users agree to data transmission)
(13) Ethical justifiability	Justifiable	Possibly restricted, because app users usually only have the choice between "agree to data transfer" or "do not use the app"

time, but we also recommend caution when using these data sources to measure the volume of tourism in a given area. At the same time, data from these sources can be useful for obtaining a more complete picture of tourism flows in a destination.

As for the statistical analysis, it is not very satisfying to be restricted to descriptive metrics because we did not include random samples of data. It might be worthwhile to check the inferential power of Bayes statistics by using one data source (e.g., local reference data) as a prior and assess the information gained through additional data sources. This, again, would require considerably more data than available here.

It has become clear that the existing studies in the field are not sufficient to systematically assess the quality and usefulness of various data sources for measuring tourism flows. We therefore suggest an assessment roadmap which includes:

(a) further discussion of the categorisation framework used in this study.

(b) further discussion of the assessment dimensions used in this study, and.

(c) the generation of more data to assess the quality and usefulness of various data sources (e.g. credit card data) in various settings and destinations.

Within argument (c), we acknowledge the necessity of not only reporting data from one source but of comparing and triangulating results from various sources, as has been demonstrated in this study. For practical and research applications, it is advisable to avoid relying on one single data source. The task of calibrating data from data sources in Categories A and possibly B with data from Categories C and possibly D would be a stepping stone in this roadmap.

At the same time, data from Categories C and possibly D could be enriched with data from Categories A and possibly B. Both approaches allow the combination of the strengths of Category A and B data sources (i.e., ability to quantify tourism flows and inter- and intra-regional tourism activity spaces) and Category C and D data sources (i.e., precise and complete spot measurements).

Declaration of Competing Interest

This research did not receive any specific grants from funding agencies in the public, commercial, or not-for-profit sectors. Mobile phone data was provided as part of a continuation of a previous research project (Reif and Schmücker, 2020). GPS data was obtained as part of an independent research project from the German Institute of Tourism Research (FH Westküste). We have no conflict of interest to declare.

Acknowledgements

The authors wish to thank the following persons for providing reference data from the destinations free of charge for research and teaching purposes: Frank Timpe (AmrumTouristik AöR), Olaf Raffel (Tourismus Marketing Service Büsum GmbH), Jan-Oliver Mauriczat (Amt Büsum-Wesselburen), Nils Stauch, Thies Jahn (Tourismus-Zentrale St. Peter-Ording), Claus von Hoerschelmann, Gerd Meurs-Scher (Landesbetrieb für Küstenschutz, Nationalpark und Meeresschutz Schleswig-Holstein), Frank Ketter, Mandy Lütt (Nordsee-Tourismus-Service GmbH), Beate Peters (Tourismusverein Westerhever-Poppenbüll e.V.). The authors wish to further thank Stefan Bornemann and Francesco Dighera at Meteonomiqs and Rahel Stöckli and Jonas Karlsson at Teralytics for constructive discussions.

References

Aagesen, H., Levin, A., Ojansuu, S., Redding, A., Muukkonen, P., & Järvi, O. (2020). Using Twitter data to evaluate tourism in Finland – A comparison with official statistics. In Muukkonen (Ed.), *Examples and Progress in Geodata* (pp. 3–16).

Adamiak, C., & Szyda, B. (2021). Combining conventional statistics and big data to map global tourism destinations before COVID-19. *Journal of Travel Research*. <https://doi.org/10.1177/00472875211051418>

Agarwal, V., Koch, J. V., & McNab, R. M. (2019). Differing views of lodging reality: Airdna, STR, and Airbnb. *Cornell Hospitality Quarterly*, 60(3), 193–199.

Ahas, R., Aasa, A., Roose, A., Mark, Ü., & Silm, S. (2008). Evaluating passive mobile positioning data for tourism surveys: an Estonian case study. *Tourism Management*, 29(3), 469–486.

Ahas, R., Armoogum, J., Esko, S., Ilves, M., Karus, E., Madre, J.-L., et al. (2014a). *Feasibility Study on the Use of Mobile Positioning Data for Tourism Statistics - Consolidated Report*. Luxembourg: Eurostat/Publications Office of the European Union.

Ahas, R., Armoogum, J., Esko, S., Ilves, M., Karus, E., Madre, J.-L., et al. (2014b). *Feasibility study on the use of mobile positioning data for tourism statistics: Consolidated report*. Luxembourg: Publications Office of the European Union.

Aparicio, D., Hernández Martín-Caro, M. S., García-Palomares, J. C., & Gutiérrez, J. (2021). Exploring the spatial patterns of visitor expenditure in cities using bank card transactions data. *Current Issues in Tourism*, 1–19.

Bastiaansen, M., Oosterholt, M., Mitas, O., Han, D., & Lub, X. (2020). An emotional roller coaster: electrophysiological evidence of emotional engagement during a roller-coaster ride with virtual reality add-on. *Journal of Hospitality & Tourism Research*, 46(1), 29–54, 109634802094443.

Batista e Silva, F., Marín Herrera, M. A., Rosina, K., Ribeiro Barranco, R., Freire, S., & Schiavina, M. (2018). Analysing spatiotemporal patterns of tourism in Europe at high-resolution with conventional and big data sources. *Tourism Management*, 68, 101–115.

Bauder, M., & Freytag, T. (2015). Visitor mobility in the city and the effects of travel preparation. *Tourism Geographies*, 17(5), 682–700.

Bauder, Michael (2019). Engage! A research agenda for Big Data in tourism geography. In Dieter K. Müller (Ed.), *A research agenda for tourism geographies* (pp. 149–158). Cheltenham: Elgar.

Birenboim, A., Dijst, M., Scheepers, F. E., Poelman, M. P., & Helbich, M. (2019). Wearables and location tracking technologies for mental-state sensing in outdoor environments. *The Professional Geographer*, 71(3), 449–461.

Bundesnetzagentur. (2021). *Teilnehmerentwicklung im Mobilfunk*. Retrieved October 25, 2021, from https://www.bundesnetzagentur.de/DE/Sachgebiete/Telekommunikation/Unternehmen_Institutionen/Marktbeobachtung/Mobilfunkteilnehmer/artikel.html?nn=268232.

Chen, J., Becken, S., & Stantic, B. (2021). Using Weibo to track global mobility of Chinese visitors. *Annals of Tourism Research*, 89, Article 103078.

Cohen, J. (1968). Multiple regression as a general data-analytic system. *Psychological Bulletin*, 70(6), 426–443.

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2. ed.). Hillsdale, NJ: Erlbaum.

Demunter, C. (2017). Tourism statistics: Early adopters of big data?. Luxembourg, from <https://ec.europa.eu/eurostat/documents/3888793/8234206/KS-TC-17-004-EN-N.pdf/a691f7db-d0c8-4832-ae01-4c3e38067c54>.

Ferber, R. (1956). Are correlations any guide to predictive value? *Applied Statistics*, 5(2), 113.

Ferreira, D., & Vale, M. (2020). Geography in the big data age: An overview of the historical resonance of current debates. *Geographical Review*, 1–17. <https://doi.org/10.1080/00167428.2020.1832424>

Gibbs, C., Guttentag, D., Gretzel, U., Yao, L., & Morton, J. (2018). Use of dynamic pricing strategies by Airbnb hosts. *International Journal of Contemporary Hospitality Management*, 30(1), 2–20.

Girardin, F., Calabrese, F., Fiore, F. D., Ratti, C., & Blat, J. (2008). Digital footprinting: Uncovering tourists with user-generated content. *IEEE Pervasive Computing*, 7(4), 36–43.

Gretzel, U. (2018). From smart destinations to smart tourism regions. *Journal of Regional Research*, 42, 171–184.

Hardy, A. (2020). *Tracking tourists: Movement and mobility*. Wolvercote Oxford: Goodfellow Publishers Ltd.

Hawelka, B., Sitko, I., Beinart, E., Sobolevsky, S., Kazakopoulos, P., & Ratti, C. (2014). Geo-located Twitter as proxy for global mobility patterns. *Cartography and Geographic Information Science*, 41(3), 260–271.

Höpken, W., Eberle, T., Fuchs, M., & Lexhagen, M. (2021). Improving tourist arrival prediction: A big data and artificial neural network approach. *Journal of Travel Research*, 60(5), 998–1017.

Izquierdo Valverde, M., Prado Mascuñano, J., & Velasco Gimeno, M. (2016). Same-day visitors crossing borders: A big data approach using traffic control cameras. In *Vortrag im Rahmen der 14th Global Forum on Tourism Statistics. Venedig*.

Janevski, T. (2019). *QoS for fixed and mobile ultra-broadband*. Hoboken, NJ: Wiley.

Kitchin, R. (2014). *The Data Revolution: Big Data, Open Data, Data Infrastructures and Their Consequences*. SAGE Publications.

van der Laken, P. (2021). ppsr: Predictive Power Score. from <https://cran.r-project.org/web/packages/ppsr/>.

Laney, D. (2001). 3D data management: Controlling data volume, velocity, and variety. *Application Delivery Strategies*. <http://blogs.gartner.com/doug-laney/files/2012/01/a949-3D-Data-Management-Controlling-Data-Volume-Velocity-and-Variety.pdf>. (Accessed 22 June 2022).

Lau, G., & McKercher, B. (2006). Understanding tourist movement patterns in a destination: A GIS approach. *Tourism and Hospitality Research*, 7(1), 39–49.

Leick, B., Kivedal, B. K., Eklund, M. A., & Vinogradov, E. (2021). Exploring the relationship between Airbnb and traditional accommodation for regional variations of tourism markets. *Tourism Economics*, Article 135481662199017. <https://doi.org/10.1177/1354816621990173>

Li, C., Zheng, W., & Ge, P. (2022). Tourism demand forecasting with spatiotemporal features. *Annals of Tourism Research*, 94, Article 103384. <https://doi.org/10.1016/j.annals.2022.103384>

Li, J., Xu, L., Tang, L., Wang, S., & Li, L. (2018). Big data in tourism research: A literature review. *Tourism Management*, 68, 301–323.

Mariani, M. M., & Baggio, R. (2021). Big data and analytics in hospitality and tourism: A systematic literature review. *International Journal of Contemporary Hospitality Management*. <https://doi.org/10.1108/IJCHM-03-2021-0301>

Mazanec, J. A. (2020). Hidden theorizing in big data analytics: With a reference to tourism design research. *Annals of Tourism Research*, 83, Article 102931.

Minkowski, H. (1910). *Geometrie der Zahlen*. Leipzig, Berlin: Teubner.

Möhring, M., Keller, B., Schmidt, R., & Dacko, S. (2021). Google popular times: Towards a better understanding of tourist customer patronage behavior. *Tourism Review*, 76(3), 533–569.

Nilbe, K., Ahas, R., & Silm, S. (2014). Evaluating the travel distances of event visitors and regular visitors using mobile positioning data: The case of Estonia. *Journal of Urban Technology*, 21(2), 91–107.

Nyns, S., & Schmitz, S. (2022). Using mobile data to evaluate unobserved tourist overnight stays. *Tourism Management*, 89, Article 104453.

O'Leary, D. E., & Storey, V. C. (2022). Data exhaust. In L. A. Schintler, & C. L. McNeely (Hrsg.), *Encyclopedia of Big Data*. Springer International Publishing.

Önder, I., Koerbitz, W., & Hubmann-Haidvogel, A. (2016). Tracing tourists by their digital footprints. *Journal of Travel Research*, 55(5), 566–573.

Padrón-Ávila, H., & Hernández-Martín, R. (2021). How can researchers track tourists? A bibliometric content analysis of tourist tracking techniques. *European Journal of Tourism Research*, 26(2601), 1–30.

- Park, S. (2021). *Big Data in Smart Tourism: A Perspective Article*, 1(3), 3–5.
- Paül i Agustí, D. (2020). Mapping tourist hot spots in African cities based on Instagram images. *International Journal of Tourism Research*, 22(5), 617–626.
- Ramos, V., Yamaka, W., Alorda, B., & Sriboonchitta, S. (2021). High-frequency forecasting from mobile devices' bigdata: An application to tourism destinations' crowdedness. *International Journal of Contemporary Hospitality Management*, 33(6), 1977–2000.
- Raun, J., Ahas, R., & Tiru, M. (2016). Measuring tourism destinations using mobile tracking data. *Tourism Management*, 57, 202–212.
- Reif, J. (2021). *Die digitale Neu-Vermessung touristischer Aktionsräume*. Dissertation. Bonn: Rheinische Friedrich-Wilhelms-Universität Bonn.
- Reif, J., & Schmücker, D. (2020). Exploring new ways of visitor tracking using big data sources: Opportunities and limits of passive mobile data for tourism. *Journal of Destination Marketing & Management*, 18, Article 100481.
- Reif, J., & Schmücker, D. (2021). Understanding tourist's emotions in time and space: Combining GPS-tracking and biosensing to detect spatial points of emotion. *Journal of Spatial and Organizational Dynamics*, 9(4).
- Revelle, W. (2021). *Psych: Procedures for psychological, psychometric, and personality research*. Evanston, Illinois, from Northwestern University. <https://cran.r-project.org/package=psych>.
- Romero Palop, J. d. D., Murillo Arias, J., Bodas-Sagi, D. J., & Valero Lapaz, H. (2019). Determining the usual environment of cardholders as a key factor to measure the evolution of domestic tourism. *Information Technology & Tourism*, 21(1), 23–43.
- Salas-Olmedo, M. H., Moya-Gómez, B., García-Palomares, J. C., & Gutiérrez, J. (2018). Tourists' digital footprint in cities: Comparing Big Data sources. *Tourism Management*, 66, 13–25.
- Saluveer, E., Raun, J., Tiru, M., Altin, L., Kroon, J., Snitsarenko, T., et al. (2020). Methodological framework for producing national tourism statistics from mobile positioning data. *Annals of Tourism Research*, 81, Article 102895.
- Sauter, M. (2018). *Grundkurs mobile Kommunikationssysteme: LTE-Advanced Pro, UMTS, HSPA, GSM, GPRS, Wireless LAN und Bluetooth* (7th). Wiesbaden: Springer.
- Schaefer, M., & Pearson, A. (2021). Accuracy and precision of GNSS in the field. In G. P. Petropoulos, & P. K. Srivastava (Eds.), *GPS and GNSS Technology in Geosciences* (pp. 393–412). Elsevier.
- Schmücker, D., & Reif, J. (2021). The big data illusion. *Zeitschrift für Tourismuswissenschaft*, 13(2), 157–166.
- Scuttari, A. (2021). Tourism experiences in motion. Mobile, visual and psychophysiological methods to capture tourists "on the move". *Tourism Management Perspectives*, 38, Article 100825.
- Shafiee, S., Rajabzadeh Ghatari, A., Hasanzadeh, A., & Jahanyan, S. (2021). Smart tourism destinations: A systematic review. *Tourism Review*, 76(3), 505–528.
- Shoval, N. (2018). Sensing tourists: Geoinformatics and the future of tourism geography research. *Tourism Geographies*, 20(5), 910–912.
- Shoval, N., & Ahas, R. (2016). The use of tracking technologies in tourism research: The first decade. *Tourism Geographies*, 18(5), 587–606.
- Shoval, N., Schvimer, Y., & Tamir, M. (2018). Tracking technologies and urban analysis: Adding the emotional dimension. *Cities*, 72, 34–42.
- Tokarchuk, O., Barr, J. C., & Cozzio, C. (2021). Estimating destination carrying capacity: The big data approach. *Travel and Tourism Research Association: Advancing Tourism Research Globally*, 51. from https://scholarworks.umass.edu/ttra/2021/research_papers/51.
- Tokarchuk, O., Gabriele, R., & Maurer, O. (2021). Estimating tourism social carrying capacity. *Annals of Tourism Research*, 86(3), Article 102971.
- United Nations Department of Economic and Social Affairs, & Statistics Division (UNSD) and United Nations World Tourism Organization (UNWTO). (2010). *International Recommendations for Tourism Statistics 2008*. New York.
- Versichele, M., Neutens, T., Delafontaine, M., & van de Weghe, N. (2012). The use of Bluetooth for analysing spatiotemporal dynamics of human movement at mass events: A case study of the Ghent Festivities. *Applied Geography*, 32(2), 208–220.
- Wetschoreck, F. (2020). RIP correlation. Introducing the Predictive Power Score, from 8080 Labs. <https://towardsdatascience.com/rip-correlation-introducing-the-predictive-power-score-3d90808b9598>.
- Witten, I. H., Frank, E., & Hall, M. A. (2011). *Data Mining: Practical Machine Learning Tools and Techniques (4. Aufl.)*. Morgan Kaufmann Series in Data Management Systems. s.l. Elsevier Reference Monographs.
- Yamamoto, M., Sato, M., & Kamitani, T. (2021). Examining spatial movement patterns of travelers: Cases in tourist destinations. In F. P. García Márquez, & B. Lev (Eds.), *Internet of things. Cases and studies* (pp. 251–273). Cham: Springer.
- Yang, X., Pan, B., Evans, J. A., & Lv, B. (2015). Forecasting Chinese tourist volume with search engine data. *Tourism Management*, 46, 386–397. <https://doi.org/10.1016/j.tourman.2014.07.019>
- Yang, Y., Fan, Y., Jiang, L., & Liu, X. (2022). Search query and tourism forecasting during the pandemic: When and where can digital footprints be helpful as predictors? *Annals of Tourism Research*, 93, Article 103365. <https://doi.org/10.1016/j.annals.2022.103365>
- Zheng, W., Li, M., Lin, Z., & Zhang, Y. (2022). Leveraging tourist trajectory data for effective destination planning and management: A new heuristic approach. *Tourism Management*, 89, Article 104437. <https://doi.org/10.1016/j.tourman.2021.104437>
- Zoltan, J., & McKercher, B. (2015). Analysing intra-destination movements and activity participation of tourists through destination card consumption. *Tourism Geographies*, 17(1), 19–35.

Dirk Schmücker is Professor for Tourism. His research interests are consumer research in tourism and empirical tourism research.

Julian Reif is a tourism researcher. His research interests are tourists spatio-temporal behaviour, urban tourism and effects of tourism.