



TECHNISCHE
UNIVERSITÄT
WIEN
Vienna University of Technology

MASTERARBEIT

Tourists vs. Locals: Mapping Urban Traces from Social Media

Ausgeführt am Department für
Geodäsie und Geoinformation
der Technischen Universität Wien

unter der Anleitung von
Francisco Porras-Bernardez, MSc., TU Wien
und
Prof. Dr. Nico Van de Weghe, Ghent University
Univ.Prof. Mag.rer.nat. Dr.rer.nat. Georg Gartner, TU Wien

durch
Yingwen Deng
Einsiedlergasse 23, 1050 Wien

03.09.2019

Unterschrift (Student)



TECHNISCHE
UNIVERSITÄT
WIEN
Vienna University of Technology

MASTER'S THESIS

**Tourists vs. Locals:
Mapping Urban Traces from Social Media**

Conducted at the Department of
Geodesy and Geoinformation
Technical University Vienna

under the supervision of
MSc. Francisco Porras-Bernardez, TU Wien
and
Prof. Dr. Nico Van de Weghe, Ghent University
Univ.Prof. Mag.rer.nat. Dr.rer.nat. Georg Gartner, TU Wien

by
Yingwen Deng
Einsiedlergasse 23, 1050 Wien

03.09.2019

Signature (Student)

ACKNOWLEDGMENTS

Here, I would like to address my greatest thanks to the support of my first supervisor MSc. Francisco Porras-Bernardez from TU Wien, my second supervisor Prof. Dr. Nico Van de Weghe from Ghent University and supervisor Univ. Prof. Mag.rer.nat. Dr.rer.nat. Georg Gartner from TU Wien. They have been offering guidance, useful remarks and encouragement through each stage of the process. I am more than grateful for the regular discussions and meetings we had. It is so impressive how supportive my supervisors have been.

Also, I would like to thank the thesis coordinator Dr. C.P.J.M Corné van Elzakker from the University of Twente. He has offered me really valuable comments and reading material which has inspired me a lot.

I would also like to acknowledge Drs. R.A. Richard Knippers from the University of Twente as the external reviewer of my thesis research.

I would also like to thank the study colleagues in Vienna. The mutual support we offered for each other means a lot.

Finally, I must express my very profound gratitude to my parents and my boyfriend for providing me with unfailing support and continuous encouragement.

Yingwen Deng,
Vienna, Austria
August 2019

ABSTRACT

Nowadays, the traces of humans are left all over the social media. Massive volunteered geographic information (VGI) contributed by social media users offers great opportunities to create new methods of understanding human activities. The generation and analysis of the digital footprints on social media have the potential to uncover the interesting spatial-temporal patterns of how people interact with the outer environment differently. The study provides an approach to differentiate the urban traces left by tourists from diverse origin countries and local citizens as different user groups based on the VGI obtained from a social media platform Flickr. Kernel density estimation was used to analyze the distribution of Flickr photos. As a result, it has been proved that it is possible to map the urban traces of tourists and locals from social media data. In addition, the approach is useful to deduce the spatial-temporal characteristics, vague region definition, and the thematic interests of local citizens and tourists from different origins.

Keywords

Social media; Volunteered geographic information (VGI); Digital footprints; Flickr; Tourists and locals; Urban traces; Kernel density estimation

Table of Contents

ACKNOWLEDGMENTS	iii
ABSTRACT	iv
LIST OF FIGURES	vi
LIST OF TABLES	vii
1. Introduction.....	1
1.1 Background.....	1
1.2 Research identification.....	2
1.2.1 Research Objectives.....	2
1.2.2 Research Questions	2
2. Theoretical background and related work.....	2
2.1 Volunteered geographic information and digital footprints	3
2.2 AOI and vague concept of places.....	3
2.3 VGI based studies	4
3. Methodology	4
3.1 Study Area and Dataset.....	4
3.1.1 Study Area.....	4
3.1.2 Original Data Review	5
3.2 Research Design	6
3.3 Overview of approach.....	7
3.4 Data Pre-processing	8
3.4.1 Flickr User Classification	8
3.4.2 Flickr Data Cleaning	11
3.4.3 Thematic POIs Filtering and Digitization	13
3.5 Data Analyzing	14
3.5.1 Approach to obtain footprints – Kernel Density Estimation	14
3.5.2 Approach to modeling city center	15
3.5.3 Approach to obtaining tourist profile	15
3.6 Data Visualization.....	17
4. Results.....	18
4.1 Footprints.....	18
4.2 Modelled City Center.....	22
4.3 Tourist Profile	24
4.3.1 Thematic interests	24
4.3.2 Temporal trend	25
5. Discussions	26
6. Conclusion	28
REFERENCE	31
APPENDIX.....	34

LIST OF FIGURES

Figure 1	Study design	7
Figure 2	Approach overview	8
Figure 3	Workflow to classify Flickr pictures.....	9
Figure 4	Workflow to obtain temporal parameters	10
Figure 5	Flickr data pre-processing	12
Figure 6	Workflow to obtain modeled city center.....	15
Figure 7	Workflow to obtain wind rose	16
Figure 8	Workflow to obtain heat map	17
Figure 9a	Footprint of the group of all tourists.....	19
Figure 9b	Footprint of the group of locals	19
Figure 9c	Footprint of the group of domestic tourists.....	20
Figure 9d	Footprint of the group of tourists from Germany	21
Figure 9e	Footprint of the group of tourists from the US	21
Figure 9f	Footprint of the group of tourists from the UK.....	22
Figure 9g	Footprint of the group of tourists from Italy	22
Figure 10	Modeled city center	23
Figure 11	Zoomed city center hotspots	24
Figure 12	Tourists profile - overview.....	25
Figure 13	The number of pictures in each season.....	26
Figure 14	Heat map.....	36

LIST OF TABLES

Table 1	Original data review	6
Table 2	Temporal parameters for different types of localness	11
Table 3	The number of Flickr data after each stage of pre-processing	13
Table 4	Seasons with corresponding months	16
Table 5	List of thematic POIs	34
Table 6	City center related tags in multi-languages	35

1. Introduction

1.1 Background

In the past, networks of fixed and mobile sensors were used to monitor and capture measurements of human's living environment and record their spatial behavior throughout the day. However, due to the prevalence of social media and the easy accessibility of the internet, the sensor network is no longer limited to those traditional ones. Humans are now involved in the sensor network. Massive geolocated data is contributed by social media users. They share their experience of how they interact with the urban environment in the form of text, image, audio or even video on various social media platforms like Twitter, Instagram, Flickr or Facebook. These data reveal the urban traces of all kinds of human activities. It offers great opportunities to create new methods of observing the environment and improving the understanding of human's spatial traces. The generation and analysis of these digital footprints can provide insights into diverse aspects such as mobility and tourism (Girardin, Vaccari, Gerber, Biderman, & Ratti, 2009).

Also, the diversity of user groups are represented on social media. Various user groups can be categorized depending on different criteria by analyzing social media data. However, in the urban environment, one of the most space-related differentiation exists between local citizens and tourists. It will definitely result in different patterns. A huge potential of such analysis can be foreseen, as these urban-specific patterns have valuable implications for both local authorities and industries like tourism. Both local citizens and tourists could benefit from it. Tourism is now considered as a new object of attention in the process of urban planning, says Jansen-Verbeke (Jansen-Verbeke, 1992). For urban planners, such studies help them plan a better urban environment with the consideration of the preservation of the environment and the harmonious coexistence of tourists and locals. For the tourism companies, for example, services like Smart Tourism Destinations can offer the right services which suit tourists' preference by the optimal usage of social media data (Buhalis & Amaranggana, 2015).

However, the core of the formation of different urban traces among local citizens and tourists is the urban areas of interest (AOIs) which refers to the areas that people are interested in. Locals and tourists have different AOIs. As a result, different urban traces are left. Closely related to the concept of AOIs, points of interest (POIs) are the points that are appealing to people. They are relevant to their visitors due to the diverse function of them. People go to certain types of POIs for certain types of services or activities. So, unlike the well-defined administrative districts, the boundary of urban AOIs are actually vague. They are regions in the mind, and it reflects how different groups perceive the environment (Montello, Friedman, & Phillips, 2014). The city center as a type of AOI is also a vague concept. It largely depends on the characteristics of the individuals or the groups which share diverse backgrounds. Different user groups like tourists from different origins and local citizens leave different urban traces on social media, and it reveals their different perceptions of AOIs like the city center.

So in this research, an approach to differentiate the urban traces left by different social media user groups based on the social media data will be the outcome. Related works will be reviewed in chapter 2. In this approach, AOIs will be extracted from the digital footprints of the local citizens and tourists from different origins. The vague concept – city center as a type of AOIs will be extracted and compared among the user groups of locals and tourists. A tourist profile regarding diverse thematic POIs will be generated. A tourist profile depicts the distinct feature of locals and tourists. In order to access this approach, it will be implemented in a case study.

1.2 Research identification

This section is to specify the research objectives of this study and the extended research questions regarding the objectives.

1.2.1 Research Objectives

The overall objective of this research is to design an approach to differentiate the urban traces left by tourists from diverse origin countries and local citizens as different social media user groups based on the VGI obtained from a social media platform.

To achieve the overall objective, it can be split into the following sub-objectives:

- a. To map the urban traces of tourists and local citizens from social media presented by their distinctive footprints
- b. To model the city center according to the semantics extracted from VGI of tourists and local citizens
- c. To create a tourist profile categorized by the origin countries of tourists as well as the local citizens in respect of the diverse thematic point of interests (POIs)

1.2.2 Research Questions

Research questions related to objective a:

Are there differences in footprints between tourists from different origins and local citizens? Which are those differences?

Research questions related to objective b:

How differently do tourists and local citizens perceive the city center?

Is there a relation between the footprints and perceived city center among tourists and local citizens? Is this relation clearer among certain user groups?

Research questions related to objective c:

Can we identify a unique tourist profile regarding different thematic POIs for different user groups?

Are there correlations between the thematic POIs in the diverse footprints and specific origin countries? Is there a seasonal trend among them?

2. Theoretical background and related work

2.1 Volunteered geographic information and digital footprints

The uploaded posts on social media with geolocated information are considered as volunteered geographic information (VGI). VGI refers to the geographic information generated and voluntarily contributed by mostly untrained private citizens who are often without qualifications. It can be considered as effective use of a sensor network which is composed of humans (Goodchild, 2007). Supported by Web 2.0 technologies, VGI as one of the most important types of user-generated web data has been a new phenomenon (Sui, Elwood, & Goodchild, 2012). A lot of attention has been drawn to the study of VGI. The nature and motivation of its producers have been studied (Coleman, Georgiadou, & Labonte, 2009; Dotan & Zaphiris, 2010). Undoubtedly, new dimensions and perspectives of geography studies (Jiang, 2013) and social science have been brought into light with the usage of VGI (Elwood, Goodchild, & Sui, 2012; Feick & Roche, 2013; Muki Haklay, 2013). The advantages and disadvantages of using VGI have been widely discussed. Comparing with traditionally acquired data, VGI has the advantages of low cost, fine resolution, covering wider geographic data and the abundance of the data amount (Wiersma, 2010). Despite all these advantages, the credibility of VGI has always been the main concern (Flanagin & Metzger, 2008). Some VGI data has been proved to have good quality. For example, Haklay compared the data from OpenStreetMap with authoritative data from Ordnance Survey (Mordechai Haklay, 2010) and the result shows a fair accuracy of the OpenStreetMap data. However, frameworks and approaches regarding crowd-sourcing, social and geographic aspects are provided to assure the quality of VGI (Fonte et al., 2015; Goodchild & Li, 2012).

The digital footprints are the locations where the social media posts are uploaded or the references of the posts to geographic entities (Stefanidis, Crooks, & Radzikowski, 2013). Compared with data on the traditional VGI platforms (such as OpenStreetMap), the digital footprints of users' geotagged posts are more of a type of indirect VGI. Because social media users share these geotagged posts mainly to share the content instead of the geographical information. In Grothe and Schaab's study (Grothe & Schaab, 2009), they propose automated approaches using Kernel Density Estimation and Support Vector Machines to generate footprints of Flickr data. In addition, the spatial distribution and densities in the urban environment are related to the topological, geometric and radial distances (Jiang, Ma, Yin, & Sandberg, 2016). Also, a number of studies were conducted based on digital footprints. For example, Salas-Olmedo used density maps to analyze the digital footprint of urban tourists (Salas-Olmedo, Moya-Gómez, García-Palomares, & Gutiérrez, 2018); the digital footprints were used to uncover mobility patterns of tourists (Girardin, Calabrese, Dal Fiore, Ratti, & Blat, 2008); it is also used to identify the tourists hotspots and evaluate the attractiveness of different spots in the city (García-Palomares, Gutiérrez, & Mínguez, 2015; Girardin et al., 2008).

2.2 AOI and vague concept of places

An AOI might be an area which contains several POIs, or just offer a nice view of other significant sights. The different intentions of different groups make every groups' AOIs

dissimilar. They visit certain categories of areas more often due to their distinct thematic interests. Since it is quite subjective to define the AOIs for individuals, the boundary of an AOI is always vague. It is largely dependent on people's cognition perception.

VGI leads to a better understanding of human activities and their perception of the environment since it is utilizing humans themselves as sensors. These data do not only contain geolocation information but also reveal thematic interests of users. Subjective opinions are encoded in the VGI, which makes it possible to extract the areas of interest (AOIs) based on it. A number of relevant researches have been conducted. For instance, a data-synthesis-driven method was adapted to extract the cognitive region of northern California and southern California (Gao, Janowicz, Montello, et al., 2017); Thematic regions were extracted based on spatial and platial user-generated data in order to identify how human defining the extent of places based on their cognition(McKenzie & Adams, 2017); city center (downtown) of Santa Barbara as a type of vague region was modeled in a vague spatial queries study (Montello, Goodchild, Gottsegen, & Fohl, 2017).

2.3 VGI based studies

Human activities can be analyzed through volunteered geographical information, especially in the urban environment. For example, large-scaled VGI obtained from Twitter was used to investigate the individual mobility and urban activity patterns (Hasan, Zhan, & Ukkusuri, 2013); Flickr dataset was used to rank the trajectory patterns in 12 different cities(Yin, Cao, Han, Luo, & Huang, 2011); urban functional regions were extracted from the aspect of human activities and POIs based on VGI (Gao, Janowicz, & Couclelis, 2017); Parks as a type of urban functional region was classified and the spatial-temporal pattern of people visiting parks was extracted based on Twitter data (Kovacs-Györi et al., 2018). As one important part of human activities, broader insights about tourism have also been offered by the usage of VGI. For instance, Flickr data was used to explore the spatial-temporal patterns of tourists' accommodation (Sun, Fan, Helbich, & Zipf, 2013); cluster analysis on Flickr photography data was used to study how tourists view the same place differently (Donaire, Camprubí, & Galí, 2014); In study of Popescu and Grefenstette, trip-related temporal information like visit time were deduced from Flickr data(Popescu & Grefenstette, 2009).

3. Methodology

3.1 Study Area and Dataset

3.1.1 Study Area

As the federal capital, one of nine states of Austria, Vienna, with more than 1.9 million ("Bevölkerung zu Jahres-/Quartalsanfang", 2019) inhabitants, is not only a primate city (Mark, 1939) in Austria but also one of the largest city in Europe. The city is well-known for its irreplaceable role in the aspect of culture, economy, and politics. As Hatz described, "If the city is evaluated by its historical significance, cultural heritage or the quality of life is falls in the

top rank” (Hatz, 2008). This specific image of Vienna enables it to be considered as one of the top destinations for tourists. Vienna has attracted over 7.5 million domestic and foreign tourists in 2018 (“Vienna Tourist Board: Arrivals & bednights 2018”, 2019). The diversity of the tourists’ background is also remarkable, which includes more than fifty countries and regions over Europe, America, Asia, Africa and Australia.

Unlike most of the tourism cities, Vienna is not solely charming for tourists but also for its local citizens. It is ranked as the top of the world’s most livable cities (“Vienna ranked as most liveable city in the world”, 2018). The recreational, as well as cultural services, are offered for both tourists and locals. In most cases, they are offered as a mixture of service for both groups. For example, the Museum Quarter, which closes the main tourist axis of the Habsburg Court, is also a much-visited gathering space for local young people (Kádár, 2013). It makes it even worthier of noting how different the urban traces are, in comparison with tourists’ ones, and taking a glimpse of the Viennese way of living.

In the meanwhile, the abundant categories of the leisure activities and tourist attractions in Vienna make it possible for each user group to have dissimilar thematic interests which can be interpreted as unique profiles. As a global and most livable tourism city, Vienna is undoubtedly suitable for this urban traces research.

3.1.2 Original Data Review

Flickr is a global online management and sharing application which is devoted to helping people make photos available to the people who matter to them and enabling new ways of organizing photos and video. By the year 2013, over 87 million registered members and more than 3.5 million new images uploaded daily (Adrianne, 2013). Notably, Flickr also provides official mobile apps for iOS, Android and an optimized mobile site, which means technically users could upload their pictures whenever and wherever is accessible to the Internet.

For this research, Flickr data in Vienna will be used. The original data have been collected by the Research Division Cartography of the Technical University of Vienna by means of Flickr application programming interface (API) and they were stored in a PostgreSQL database. In this dataset, attributes of both pictures and their owners are included (as shown in Table 1). Starting with the attributes of pictures, the dataset has the distinctive photo IDs, title, the created dates, related semantic tags, number of views and the relevant geo-location of pictures. For the photo owners, the dataset includes their user IDs, origins indicated in the user profiles in multiple languages, the claimed origin countries in English processed with the GeoNames geographical database (GeoNames, 2019) and the classified country of origin (in English). As it is explained in the paper (Verstockt et al., 2019), the country in which the user had uploaded most pictures for a period greater than 6 months is classified as the origin country of this user.

The dataset contains 479,126 pictures of 13,187 users from 117 countries/regions in total. The temporal duration of all the data is from January 2nd, 2002 to December 5th, 2018. Therefore, the overall interval of the used dataset is around 17 years.

When Flickr users upload their pictures, they can add a description about the picture, add semantic tags, tags of people who are in the pictures as well as geotag the pictures. When a picture is geotagged, the contents are related to a specific location on the Earth's surface (Goodchild, 2007). With the geotags of the pictures, the precious relevant locations of the pictures are recorded. However, the tagged locations are not always exactly where the pictures are taken. Only when users are uploading pictures taken by the auto-geotagging capable device, the tagged locations are always exactly where the photos are taken. But with pictures no automatically added geotags, the tags can be later added. In this case, it is possible that the tagged locations are not exactly where the pictures were taken, which is commonly regarded as the issue of the accuracy of VGI. For this study, all locations with precise coordinates are considered as the location where pictures are taken. And the accuracy of the VGI and its effects will be discussed in the later session.

Flickr Picture			Flickr User		
Attribute Name	Data Type	Meaning	Attribute Name	Data Type	Meaning
photo_id	integer	The unique identifier of pictures	photo_owner	varchar	The unique identifier of Flickr users
title	text	The title of pictures assigned by users	profile_locat	varchar	The user-defined origins for users' profile (multilingual)
date_taken	date	The date when the pictures were taken	profile_processed	text	The result of profile_locat processed by GeoNames (in English)
tags	text	The semantic tags assigned by users for each picture	country_classif	text	The result of the classification of users' origin countries (in English)
views	integer	The number of a picture being viewed by Flickr users			
point	geometry	The picture data stored as points with their geographic location			

Table 1 Original data review

3.2 Research Design

This study is designed as five main stages (shown as Figure 1), which includes original data review, data pre-processing, data analysis, result visualization and result conclusion. The data pre-processing is divided into two parts, which are the processing of Flickr Data and the processing of thematic POIs. In order to achieve the three sub-objectives mentioned above, the data analysis phase is also separated into three parts including obtaining the footprints of different user groups, modeling the city center and creating a tourist profile. As for the visualization, a density map is used to represent the footprints of tourists from different origin

countries and also local citizens; A wind rose is used as a representation of tourist profile which depicts the thematic interests of each user groups; A heat map is used to illustrate the temporal trend of the footprints among different user groups. The detailed study approach will be explained in the later section.

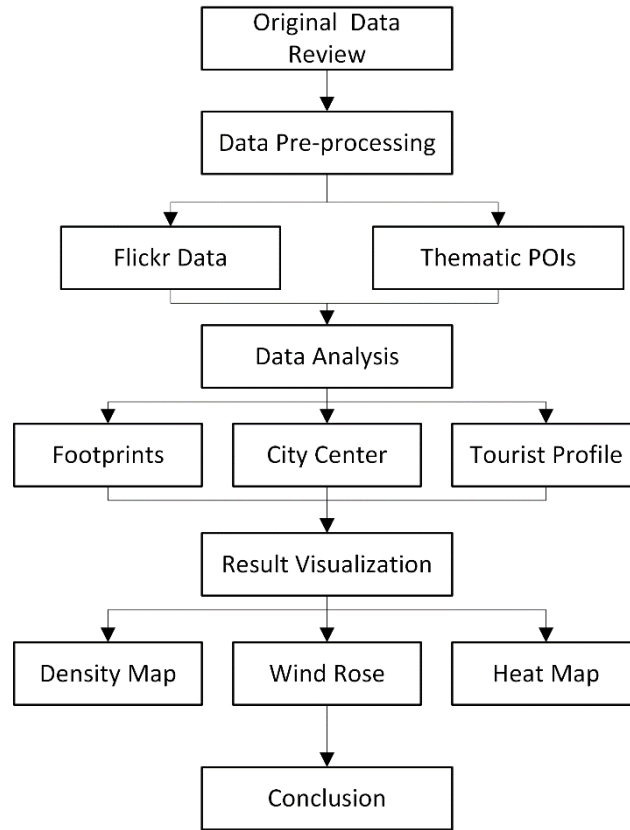


Figure 1 Study design

3.3 Overview of approach

In order to solve the research questions, the study approach (see Figure 2) is designed to obtain the footprints by applying Kernel Density Estimation on the classified Flickr data in order to represent the urban traces of different user groups.

With a threshold filtering, the areas of interest (AOIs) of each user group can be extracted. Comparing with the extracted AOIs with thematic POIs, the characteristics of each user group can be concluded into a tourist profile. In the meanwhile, the seasonal trend of each groups' urban traces can be revealed by the footprints obtained from temporal-wise classified Flickr data as a temporal aspect of the tourist profile. And city center as one type of AOI can be modeled from semantically filtered Flickr data.

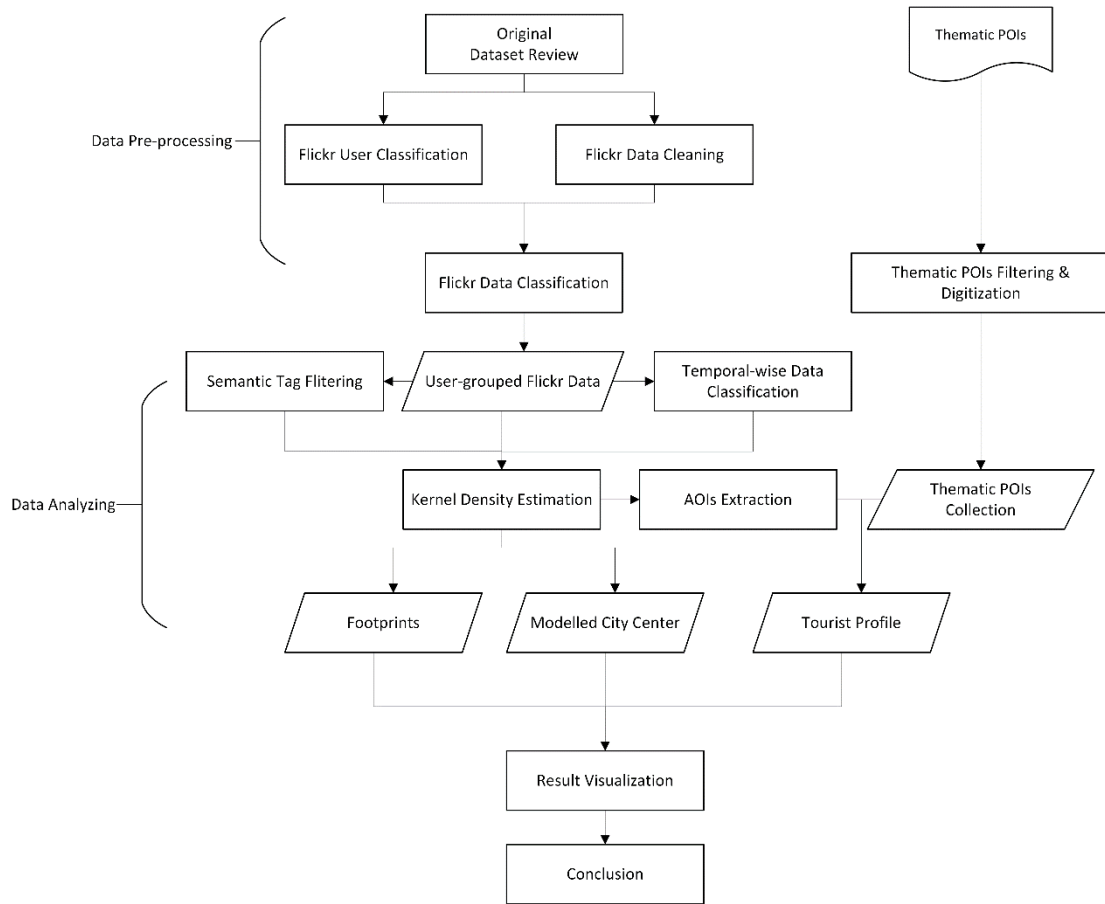


Figure 2 Approach overview

3.4 Data Pre-processing

As mentioned in the research design, the pre-processing of data is divided into two parts: One is the processing of the original Flickr dataset which includes the classification of Flickr users and the aggregation of picture data; the other is the processing of the thematic POIs collection which includes the filtering and digitization of the POIs. The detailed description of the approach will be included in the following three sections.

3.4.1 Flickr User Classification

Although the origin countries of each Flickr user are included in the original data, it is insufficient to directly classify the locals and domestic tourists. Because they are all from Austria and the precise cities are sometimes absent from the user profile. There are generally three types of Flickr users in the original data, which are those without information about classified origin countries, users from other countries and Austrian users. And since Germany, the US, the UK and Italy are the countries that contribute most tourists to Vienna, international users from these countries are considered as main study user groups in this study. Pictures uploaded by users from these countries are extracted separately. But for Austrian users, there are “known_locals” with profile indicating their origin as Vienna, “unknown_aut_owners” who are classified by the algorithm as Austrians in the original data (as mentioned in chapter 3.2.1) without indicating any precise cities and “known_aut_tourists” who have set their origins as

other cities in Austria. However, for “unknown_aut_owners”, the profiled origins could be a null value or simply indicating Austria. It is easy to directly distinguish the users from other countries, so the main purpose of this user classification step is to classify the locals and domestic tourists from the “unknown_aut_owners” depending on the featured temporal parameters extracted from pictures of “known_locals” and “known_aut_tourists” (see Figure 3). Before calculating the featured temporal parameters, users with their related information are extracted from the pictures they contributed for the purpose of extracting the classifying temporal parameters. This information includes the ID of the users, the number of distinct dates on which they uploaded pictures, the number of their pictures, the ascending time sequence of all the distinct uploading dates, the newest and oldest uploading dates as well as the calculated durations.

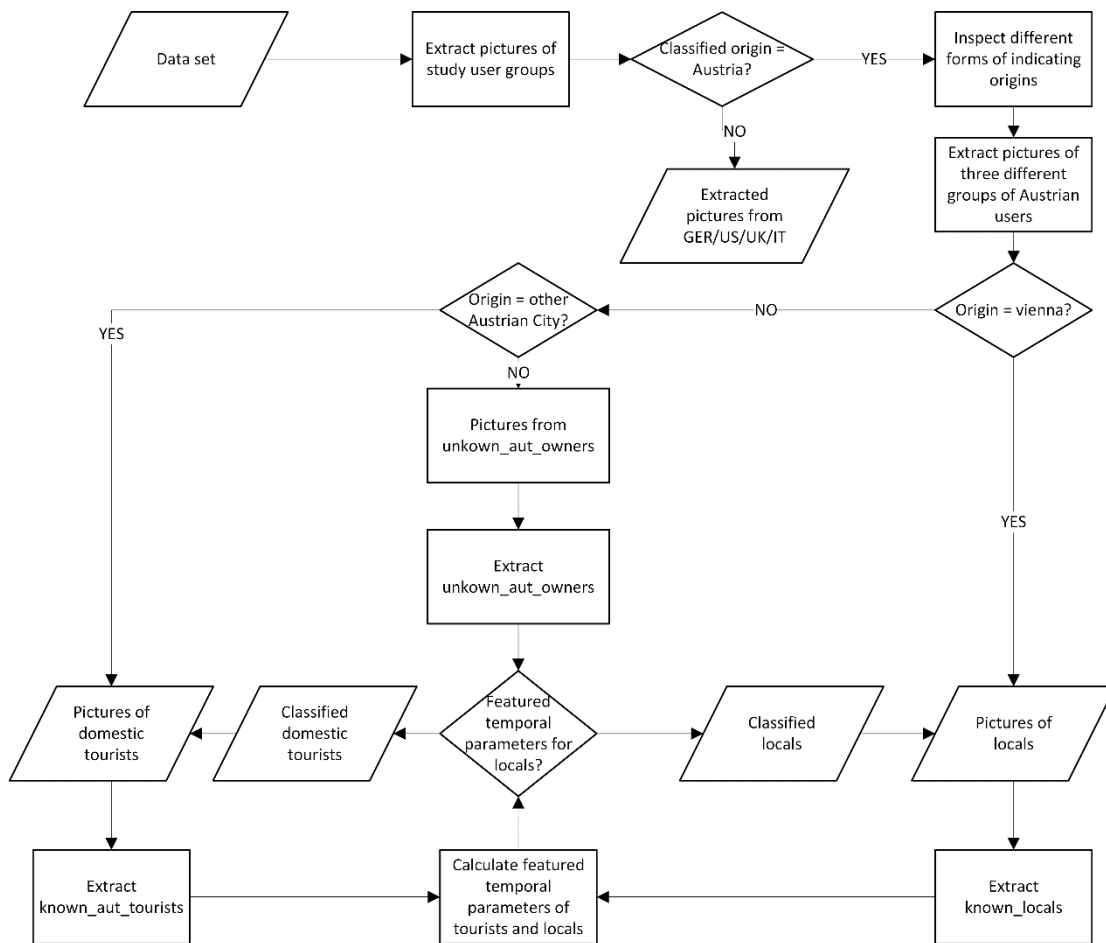


Figure 3 Workflow to classify Flickr pictures

The featured temporal parameters used to classify the unknown users are average duration, maximum intervals of users uploading pictures and their average visit time. Assuming that all Flickr users upload pictures actively during their visit to Vienna, these parameters can reveal how long these users spend their time in Vienna. The relational database management system used here is PostgreSQL. The programming language Python is used to conduct the calculation (see Figure 4).

- Average duration

To calculate the duration of a user visit Vienna, the time difference between the newest uploads and oldest uploads is calculated. The average duration is the average value of all the durations of users belonging to the same user group.

- Maximum intervals

Interval is the period of time between the two pictures that are uploaded. For each user, the maximum interval is the longest duration during which there are no uploaded pictures. However, the maximum intervals here are the average values of all the maximum intervals of all users from one user group.

- Average visit time

The visit time is a deduced value from the calculated intervals which represents the time from the user's arrival to the departure. Any intervals longer than 60 days would be considered as the periods when users leave the city. So for multiple-visit users, the date before the long interval is considered as the departure date. The first date after the long interval when the user uploads a picture is considered as the arrival date for the next visiting time. For each user, the average visit time is apparently the average of all their visit time. But the parameter – average visit time here is the average value of all average visit times of all users from one user group.

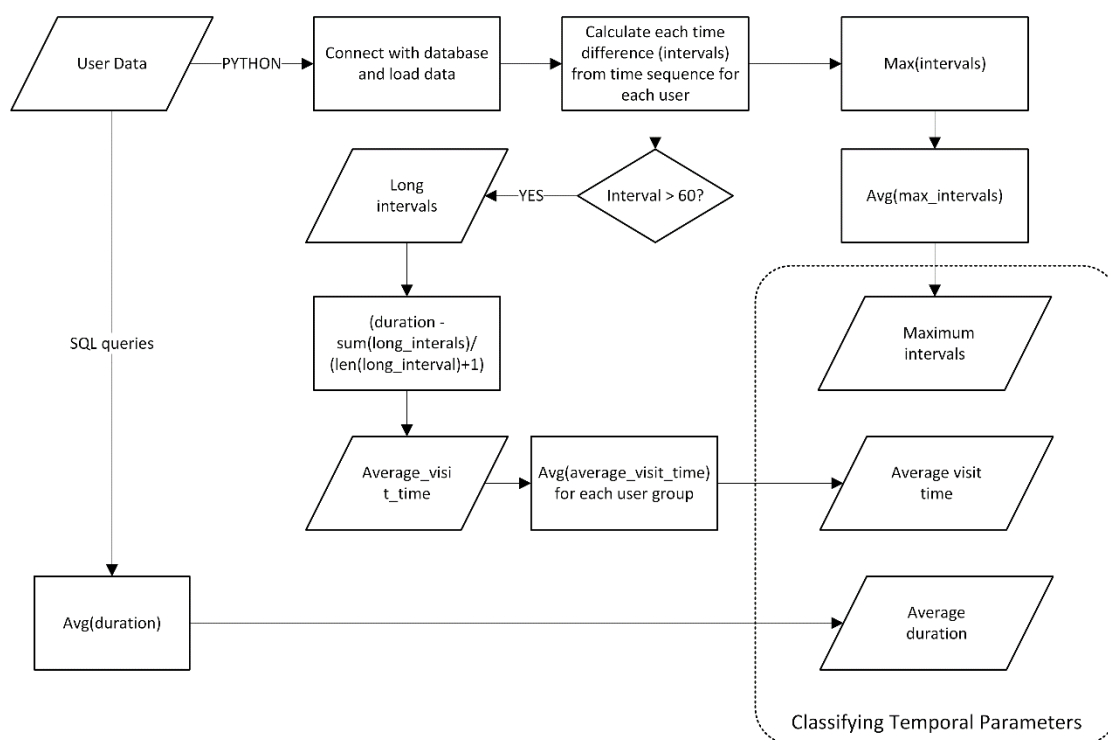


Figure 4 Workflow to obtain temporal parameters

In Zhang's research (Zhang, 2019), thresholds of temporal parameters regarding different localness types of Twitter users are defined. Compare with these thresholds, the obtained ones here differ a lot (see Table 2). However, due to the different characteristics of different social

media platforms, the threshold of classifying users differs naturally. So, depending on the obtained parameters, 36 locals from “unknown_aut_owners” are extracted with the remaining 1203 users as domestic tourists.

	Average duration	Maximum interval	Average visit time
Flickr Users			
Known_Locals from data	>1026 days	<598days	>46 days
Twitter Users			
Long-term residents	>365 days	<60 days	-
Temporary/short term residents	30-365 days	<60 days	-
Seasonal resident	>=730 days	>= 180 days	30-90 days
Non-local commuter	>= 30days	< 60 days	-
Visitor (once)	< 30days	-	-
Visitor (multiple times)	-	60-180 days	<30 days
Tourists	<= 7 days		

Table 2 Temporal parameters for different types of localness

3.4.2 Flickr Data Cleaning

The major function of this data aggregation is to prepare the data for later phases. To reduce the data bias is the main purpose of it. The main bias source lies in the data contributing behavior of social media users. As we know, the activeness of social media users differs a lot. Active users contribute a great number of pictures while less active ones might only upload a few. In this dataset, for example, an active local Flickr user uploaded more than 28,000 pictures all together while the least active local user only uploaded one picture. In spite of the considerable amount of data provided by active users, the resulting footprints will be dominated by the behaviors of these active users while the behavior of the inactive ones is overlooked (Hu et al., 2015). Two main stages of data cleaning are adapted to avoid the domination of active users.

3.4.2.1 Data Aggregation

Shown as Figure 5, the first stage of the Flickr data cleaning is the data aggregation. After reviewing the original data, it appears to be possible that some pictures uploaded by the same users are sharing the same geo-location. For example, on June 11, 2017, 936 pictures uploaded by one extremely active user are overlapped. Such a case should be avoided undoubtedly. With SQL queries, spatially overlapping pictures are grouped by users and date. Specifically, overlapping pictures uploaded by the same user on different dates are not aggregated in order to keep the temporal trend of footprints. For each group of overlapped pictures, one with the most time of view is selected as the representative which means only one picture is reserved by each group. To avoid losing semantic tags due to the aggregation, tags are aggregated for each group and then assigned to each representative. As a result, the original data set is reduced to 208,348 pictures. Austria, Germany, the United Kingdom, Italy, and the United States are five top countries contributing most visitors to Vienna.

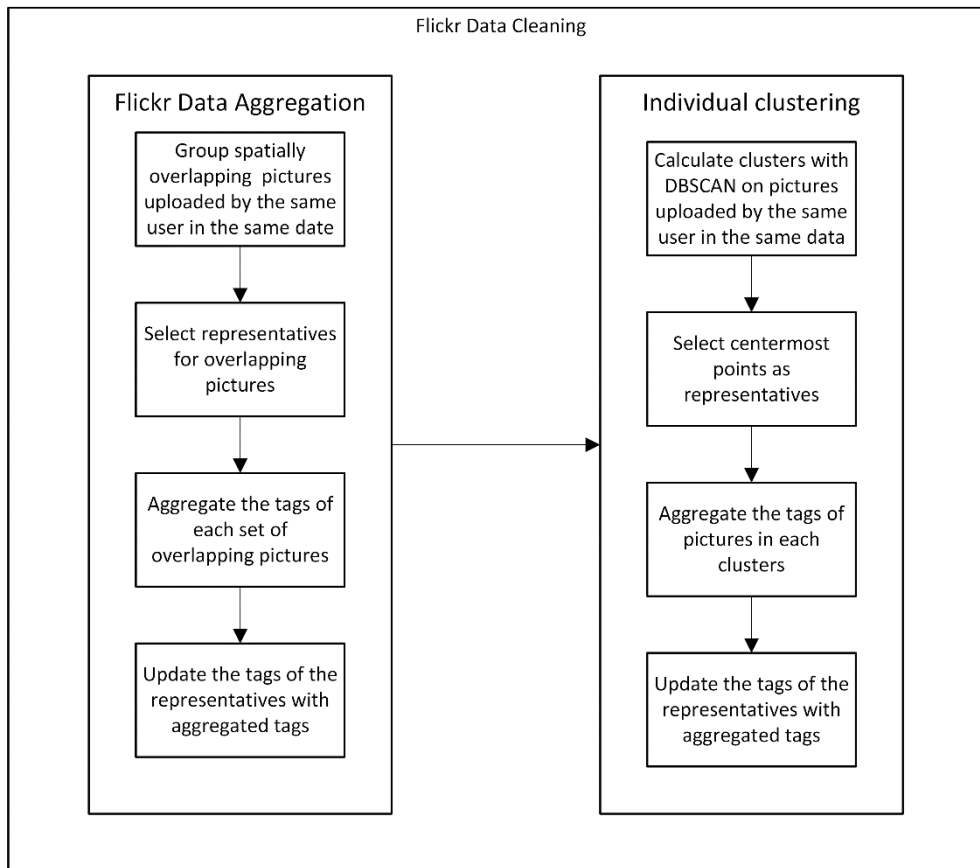


Figure 5 Flickr data pre-processing

3.4.2.2 Individual clustering with DBSCAN

The second stage is the individual clustering using DBSCAN (density-based spatial clustering for application with noise) which aims at reducing the effect of extreme contributors who upload a large number of pictures in a small radius of the area. Although these pictures are not overlapped, they still have a powerful effect on the extracted footprints. Few extreme contributors might have strong affection on some certain POI and, as a result, these few contributors would dominate the patterns of a specific local region (Gao et al., 2017). By applying DBSCAN with SQL queries, pictures uploaded by the same users concentrated in a small region can be clustered into different shapes and replaced by the centermost points as representatives. The domination of the extreme contributors is reduced as the number of uploads contributed by each user regarding each POIs is reduced. In the meanwhile, since the clustering is only applied to the pictures from the same user as well as sharing the same date, the temporal trend of footprints is preserved. After the clustering and representative selection, the tags are aggregated and updated as the former steps in the data aggregation.

Regarding the clustering method, DBSCAN is a density-based clustering method that is superior for processing spatial data (Boeing, 2018; Sun, Fan, Helbich, & Zipf, 2013). Unlike other clustering methods, DBSCAN does not require a predefined number of clusters. And in this study, the number of picture clusters is difficult to estimate since the number of each user's

uploads is different. Also, with DBSCAN, the clusters are not limited to one convex shape but in arbitrary shapes. So pictures along a pedestrian or round court regarding one POI can be clustered into a linear or circular cluster. As Kriegel depicted, “A density-based cluster is a set of data objects spread in the data space over the contiguous region of a high density of objects, separated from other density-based clusters by contiguous regions of low density of objects”(Kriegel, Kröger, Sander, & Zimek, 2011). So in this procedure, pictures distributed in a higher density region are clustered and separated with clusters of pictures with the lower density of objects instead of being clustered based on their distances to the centroid.

By applying DBSCAN, the proper value of the search radius -- Eps and the minimum number of points within the search radius -- MinPts are needed. In this case, the Eps is decided based on the geographic scale of the research area and the distribution of POIs. After several experiments, Eps is selected as 30 meters. As for the MinPts, it is assigned to 1. So every picture point is assigned to either a cluster or forms its own cluster of size one. In this way, the sole picture points in the low-density region are kept since they are not classified as noise.

To obtain the centermost point, the centroids of each cluster are calculated and matched with the points within the corresponding clusters. After calculating the distances of each point to the centroid, the points with minimum distance to each centroid are selected as the centermost point. As a result, the remaining amount of pictures from each study user group is shown in the table below (Table 3).

Origins of User groups	After aggregation	After individual clustering
Vienna (Local)	58,552	46,285
Abroad (All tourists)	149,797	116,216
Austria (domestic tourists)	34,874	28,238
Germany	12,826	10,083
The United States	12,884	9,094
The United Kingdom	11,007	8,274
Italy	8,158	6,833

Table 3 The number of Flickr data after each stage of pre-processing

3.4.3 Thematic POIs Filtering and Digitization

To obtain the thematic POIs in Vienna, the official website of the tourism board for Vienna is referenced. On this website, suggestions for sightseeing (“Sightseeing in Vienna”, 2019) and other entertainment activities like dining and drinking (“Appreciating Vienna”, 2019) are offered. As result, there are in total 8 categories including shopping areas, religious sights and architecture, operas and theaters, nature and parks, museums, historical sights and architectures, dining and drinking as well as contemporary sights and architectures. But for some categories, the POIs are too dense, so a selection depending on the overview of the distribution of those points are needed. There are altogether 64 spots (see appendix table **) on the list.

In addition, since most of the POIs are actually areas and extracting them into representative points could possibly result in information absence, the final extracted shapefile of these thematic POIs is actually a shapefile of polygons. To extract the shapefile, the newly updated orthophotos provided by the Stadt Wien website (Stadt Wien, 2019) are used as the base map and the polygons of the filtered POIs are digitized with ArcMap. The orthophotos are obtained with a 15cm resolution in the duration of March 2018 to April 2019.

3.5 Data Analyzing

3.5.1 Approach to obtain footprints – Kernel Density Estimation

The urban traces left by tourists and locals are represented by the footprints of uploaded Flickr pictures. The analysis of the point pattern - the footprint is to explain the empirical spatial distribution of Flickr data points in order to infer the underlying spatial point process, which in this case is the diverse visiting behavior of tourists and locals. There are two interrelated approaches to describe such a point pattern: the first-order effects reflect the intensity of points, while the second-order refers to the interaction between points (O'sullivan & Unwin, 2014). In this study, the location of each point of uploaded pictures is assumed to be independent and the spatial association between each point is not considered (Sun et al., 2013). So the footprints of Flickr pictures are depicted by continuous surfaces of diverse concentrations of pictures.

However, before obtaining the footprints, another aggregation for the pictures is necessary. As mentioned in the pre-processing phase, the previous aggregation of pictures only aggregate overlapped pictures owned by the same users uploaded on the same date in order to keep the temporal trend. But the footprints to be obtained here are based on all the pictures from the whole duration, so the temporal trend is not relevant here. Moreover, it is possible that there are overlapped pictures owned by the same users but uploaded on the different dates are still preserved. To avoid bias, this aggregation is needed to eliminate the excess overlapped pictures. To obtain these footprints, Kernel Density Estimation is applied with the spatial analyst tool provided by ArcMap 10.6.1. KDE is commonly used for geospatial information analysis to estimate the density distribution of the geographic process (O'sullivan & Unwin, 2014). It calculates a magnitude-per-unit area from point features using a kernel function to fit a smoothly tapered surface to each point (“Kernel Density”, 2019).The function:

$$f(x) = \sum_{i=1}^N \alpha_i k_h(x - x_i)$$

Equation 1

returns the estimated density at x . The α_i is the kernel weights with $\sum_{i=1}^N \alpha_i = 1$. Normally, all kernels are equally weighted as $\alpha_i = 1/N$. The kernel function $k_h(\cdot)$ is required to satisfy $\int k_h(x)dx = 1$ and $k_h(x) \geq 0 \forall x \in R^2$. And the h is the bandwidth parameter that determines the smoothness of the surface (Grothe & Schaab, 2009). While in this process, the bandwidth h is chosen individually for each dataset of different user groups. The adapted bandwidths are the default search radius which is calculated based on the spatial configuration and number of input points.

3.5.2 Approach to modeling city center

As a subjective fuzzy concept, the city center can be modeled by obtaining footprints of city center related pictures. And to describe the city center, KDE is conducted on the related pictures to show the distribution of the density of pictures. As a result, smooth surfaces with the diverse density of city center related pictures are acquired for tourists and locals. The areas with higher densities of pictures represent the areas that are more commonly considered as the city center. The filtering of city center related pictures involves the semantic filtering of tags. Due to the diversity of expressing city center and languages, a list of multilingual tags related to the city center is generated manually (see appendix table **). With SQL queries, pre-processed Flickr pictures containing those tags are extracted from classified pictures from all tourists and locals. Same with the footprints obtaining, the city center is also modeled on all the pictures from the whole duration, so temporal trend is not relevant here. To prevent bias, another aggregation to eliminate the excess overlapped points is also necessary before applying KDE. It is adapted after the semantic filtering. After the semantic filtering and aggregation, 898 pictures are filtered out from all the locals' uploads while 1,721 pictures from all tourists remain. At last, KDE is applied to the aggregated data of both groups (see Figure 6).

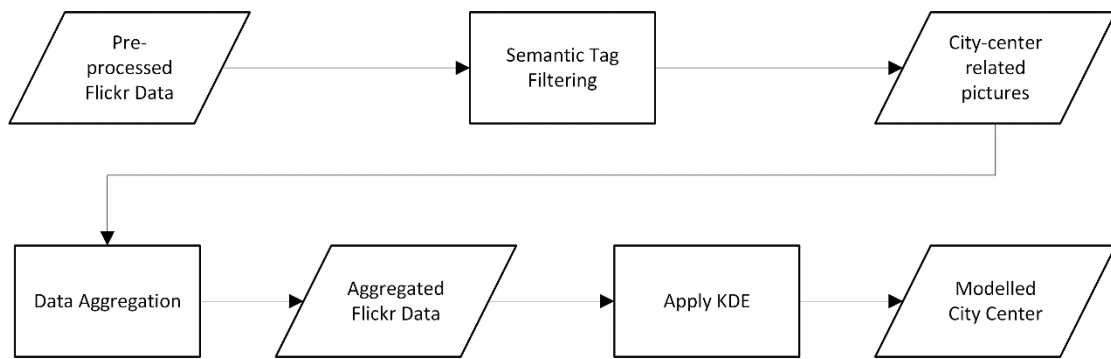


Figure 6 Workflow to obtain modeled city center

3.5.3 Approach to obtaining tourist profile

There are two aspects of the tourist profile. One is the general overview of each user group's thematic interests with is represented by the wind rose. The other one reveals the temporal trend of their footprints by the means of the heat map.

To obtain the thematic interests of different user groups, a threshold filtering is applied to the KDE result. In this case, areas with a density of pictures higher than 30% are considered as Flickr users' AOIs. Comparing the extracted AOIs of each user group with the shapefile of thematic POIs, if the polygon which represents the POI is overlapped with any area of picture's density higher than 30%, then it is considered as a targeted POI for the corresponding user group. By statistically analyzing the ratio of targeted POIs of each category, the thematic interests of each user group can be depicted by the wind rose as an aspect of the tourist profile (see Figure 7). The larger the ratio is, the more interested a user group is at a certain category of sights.

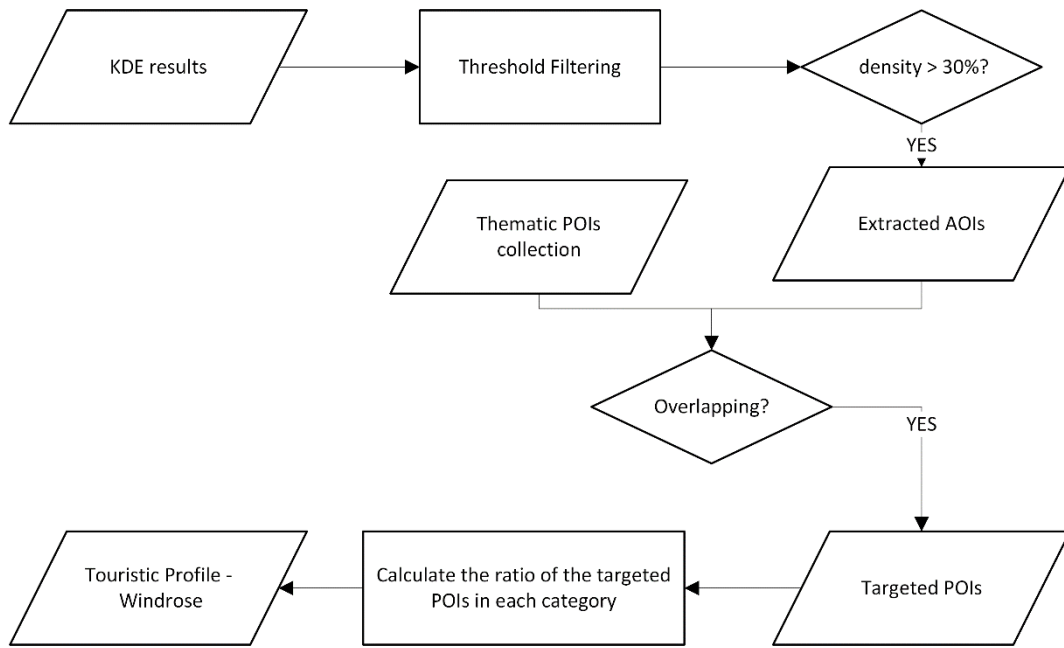


Figure 7 Workflow to obtain wind rose

Due to the available temporal information of the Flickr data, the temporal trend to be studied here is the seasonal trend. The pre-processed Flickr data is extracted and grouped into four seasons with SQL queries. Considering the climate characteristic of Austria, pictures uploaded in certain months are considered as pictures of a corresponding season (see Table 4). After grouping, KDE is applied to each group of pictures. Smooth surfaces depicting the density of uploaded pictures for each season are obtained. Comparing the obtained KDE results with the thematic POIs, we can see how dense the pictures are regarding one POI. The density is assigned as a feature value to the POI. The higher the value is, the more attractive the POI is considered.

Season	Months		
Spring	March	April	May
Summer	June	July	August
Autumn	September	October	November
Winter	December	January	February

Table 4 Seasons with corresponding months

There are occasions when one polygon is overlapped with areas of multiple densities. In this case, the maximum density is considered as the featured value. But in some cases, the POI meets with a higher density area at its boundary. To determine which density this POI belongs to, the potential location of taking pictures needs to be considered. If it is possible to take the pictures of this POI along the meeting boundary, then the value of the higher density is considered as the feature value. If not, the value of the maximum density area overlapped with the POI is the feature value (see Figure 8).

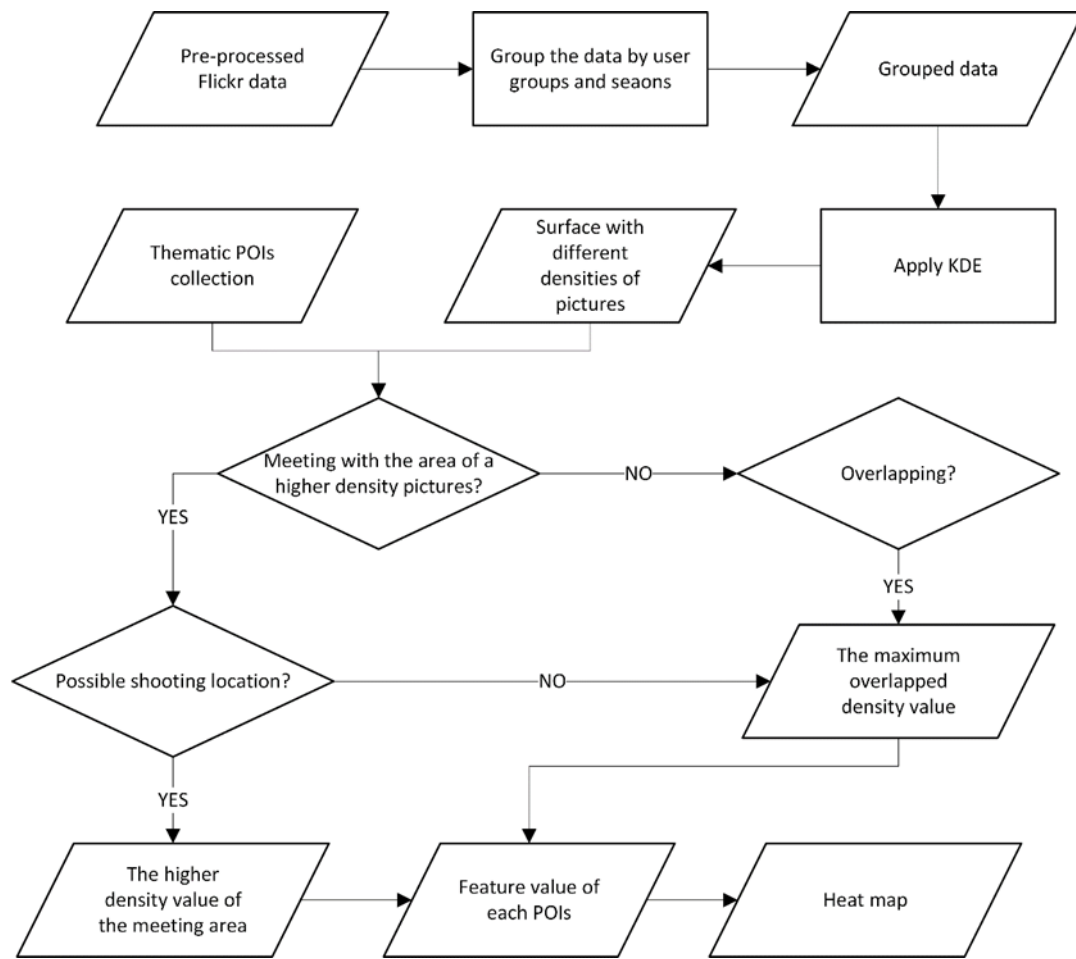


Figure 8 Workflow to obtain heat map

3.6 Data Visualization

As mentioned in the research design, density map, wind rose plot and heat map are used to visualize the analysis results. These three visualization approaches are introduced briefly in this chapter.

Density maps are used to render the density difference visually (Bertini, Di Girolamo, & Santucci, 2007). The color scale is commonly used to represent different density values on the map. Recently, it has been quite common to apply density mapping using “big data” in order to find densities of certain phenomenon (GIS Lounge, 2017). In this study, density mapping using kernel density estimations is applied to show the spatial distribution of the uploaded pictures of different groups. These acquired footprints depict the urban traces left by Flickr users. The higher picture density areas represent the areas with higher rates of Flickr users’ activities in the urban environment. The kernel density estimation tool in the spatial analyst toolbox of ArcGIS 10.6.1 is used here to generate the density maps.

A wind rose is a graphic tool used by meteorologists to summarize information about how the wind blows from each direction during the observation period (Encyclopaedia Britannica,

2013). In some wind rose graphs, the distribution of wind speed in each direction is also depicted. A longer spoke in one direction means a higher frequency of wind. However, in this study, eight directions in the wind rose plot represent eight thematic POIs categories. And a longer spoke means a higher ratio of targeted POIs in that category for certain user groups. Microsoft Excel is used here to generate the wind rose plot.

The heat map is a way to visualize data by representing the individual values contained in a matrix as colors. Larger values are represented by squares with darker colors. A Python data visualization library based on matplotlib, Seaborn is used to plot the heat map (“Seaborn.Heatmap”, 2019). In this study, the feature values of all POIs in different seasons are represented as colors.

4. Results

4.1 Footprints

The footprints of each user groups’ uploads depict their unique urban traces. The KDE results of each user group are calculated with default bandwidth (see chapter 3.3). And obviously, for each user group, some hotspots with a higher concentration of pictures can be noticed.

Among all the obtained tourists’ footprints, a similar overall pattern is shared: Despite the different density value distribution, it is noticeable that pictures are always concentrated in the southwest part of Innere Stadt, Schönbrunn, Belvedere, and the northwest corner of Prater (see Figure 9a). Especially, the area around Stephansdom is overlapped with the area where more than 90% pictures are located in all footprints; areas around Imperial Palace and Heldenplatz (Heros’ Square) are always overlapped with the area of more than 40% picture density in all footprints as well (including locals’ footprints).

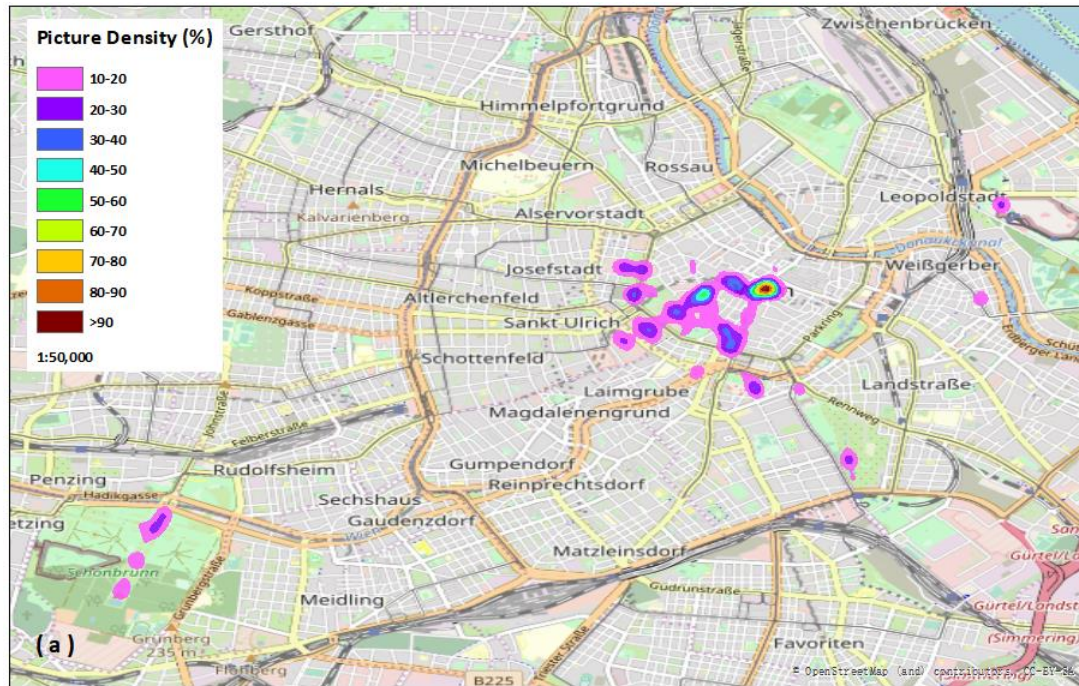


Figure 9a Footprint of the group of all tourists

As for the footprints of local citizens (see Figure 9b), they are more dispersed than other footprints. More areas with a relatively higher density of pictures are revealed. The boundary of areas with more than 20% picture density is expanded to the district Liesing as well as the district Kagran. For example, regions around Kaisermühlen and Simmering are both one of the newly revealed hotspots. Also, for local citizens, pictures are more concentrated in the southwest part of the district Innere Stadt. The area of picture density greater than 40% is larger, compared with the area with the same picture density in the footprints of other user groups.

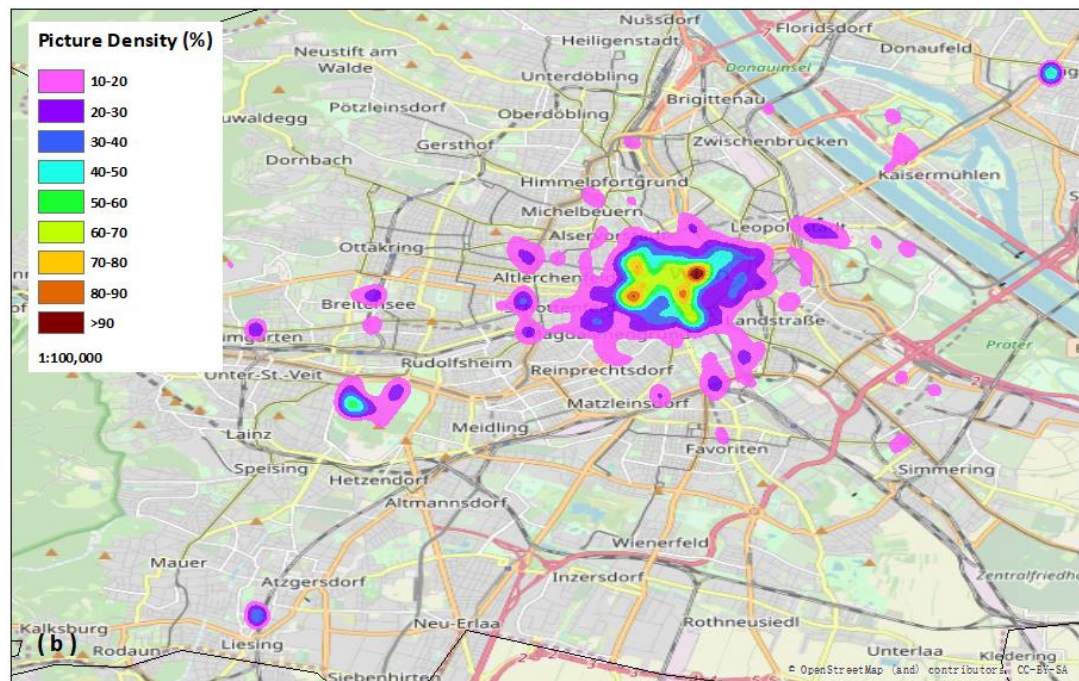


Figure 10b Footprint of the group of locals

Regarding the footprints of domestic tourists (Figure 9c), it also shares certain similarities with the locals' footprints apart from the footprints of other tourist groups. Hotspots in the district Kagran, Kaisermühlen can be also noticed. Same with the locals' footprint, there is a higher concentration of pictures at Schönbrunn zoo while you can barely see it in footprints of other tourist user groups. But for the area at the northwest corner of Prater, either the footprints of domestic users nor local citizens show a higher density of pictures comparing with other user groups.

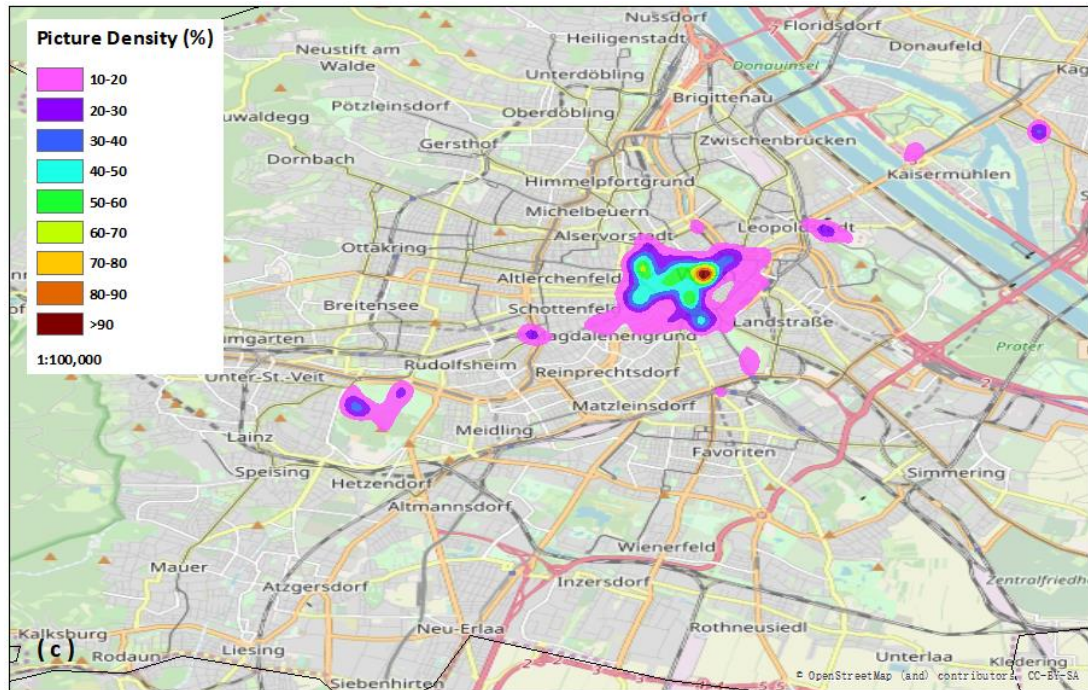


Figure 11c Footprint of the group of domestic tourists

For the footprints of other tourist groups, the patterns are generally similar. It is worth noticing that pictures are denser (greater than 20% of density) in the area of Belvedere in the footprints of tourists from Germany, the UK as well as Italy than the groups of US and Austria (in the interval of 10% -20% picture density). In addition, it shows a higher concentration of pictures at the northwest corner of Prater in the footprints of tourists from Germany and the UK. On the contrary, the footprint of tourists from the US shows a relatively lower concentration of pictures (Figure 9d, 9e, 9f, 9g).

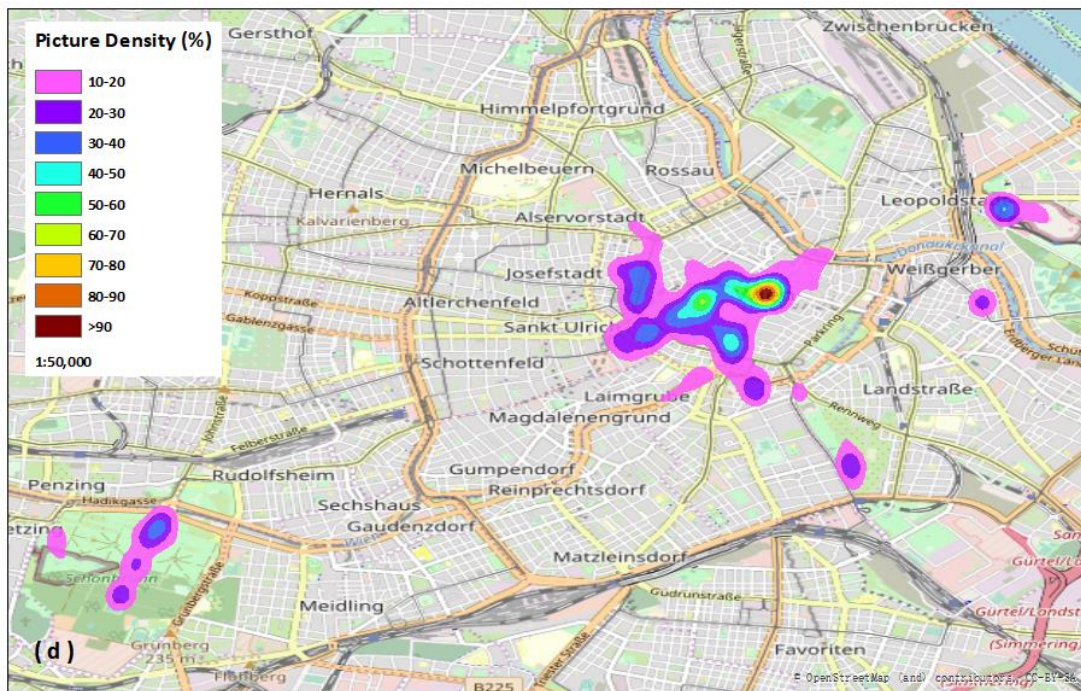


Figure 12d Footprint of the group of tourists from Germany

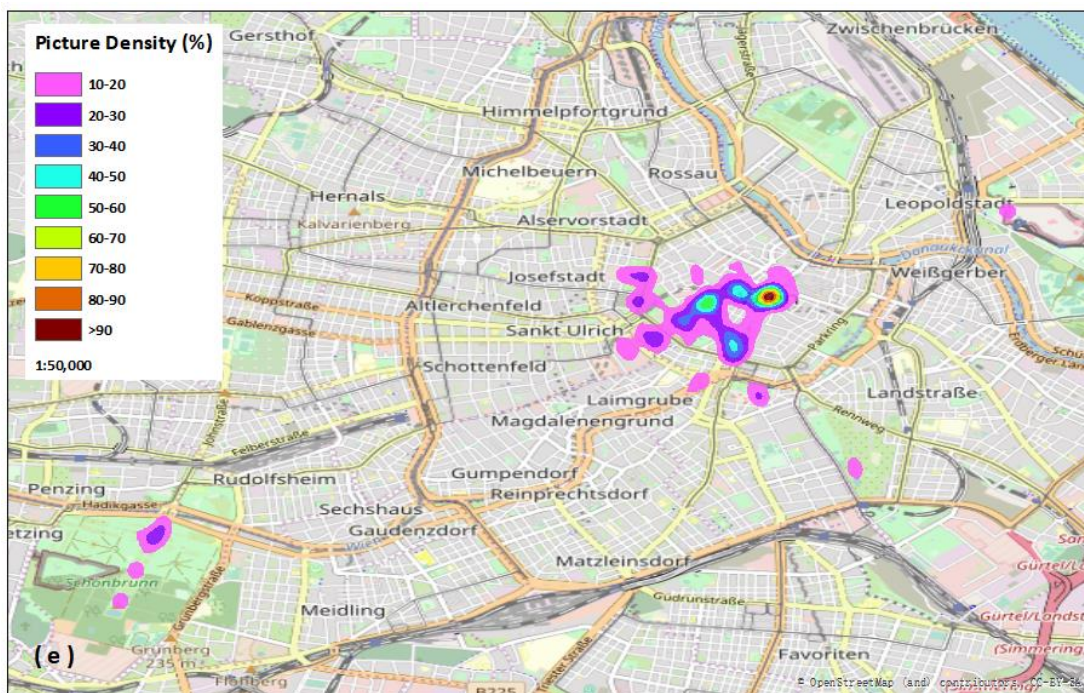


Figure 13e Footprint of the group of tourists from the US

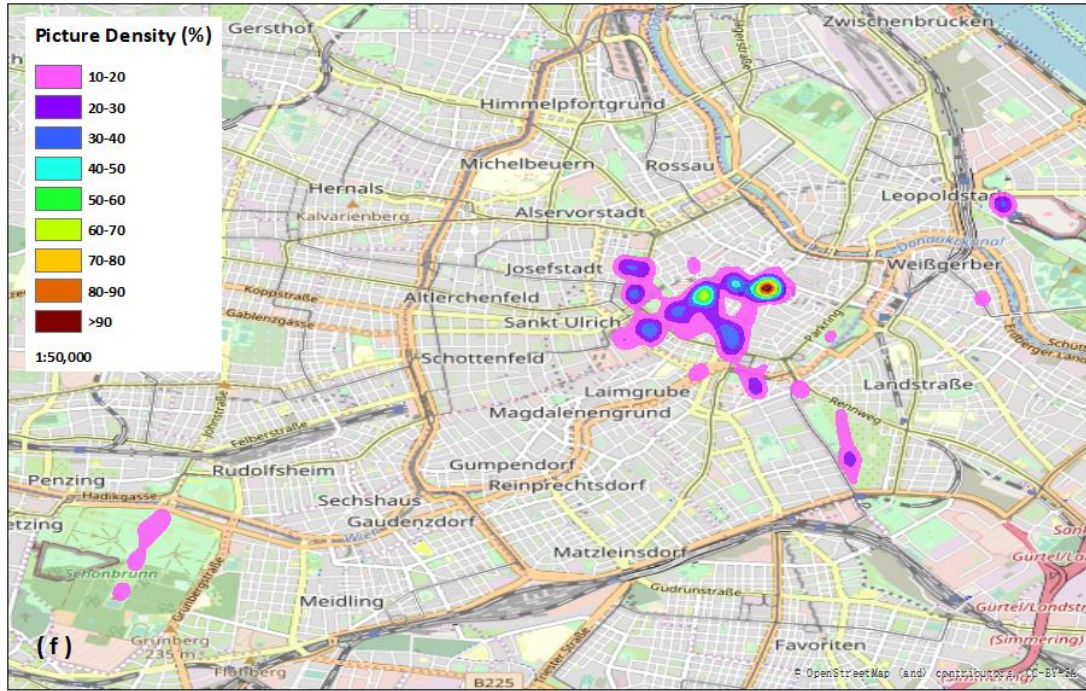


Figure 14f Footprint of the group of tourists from the UK

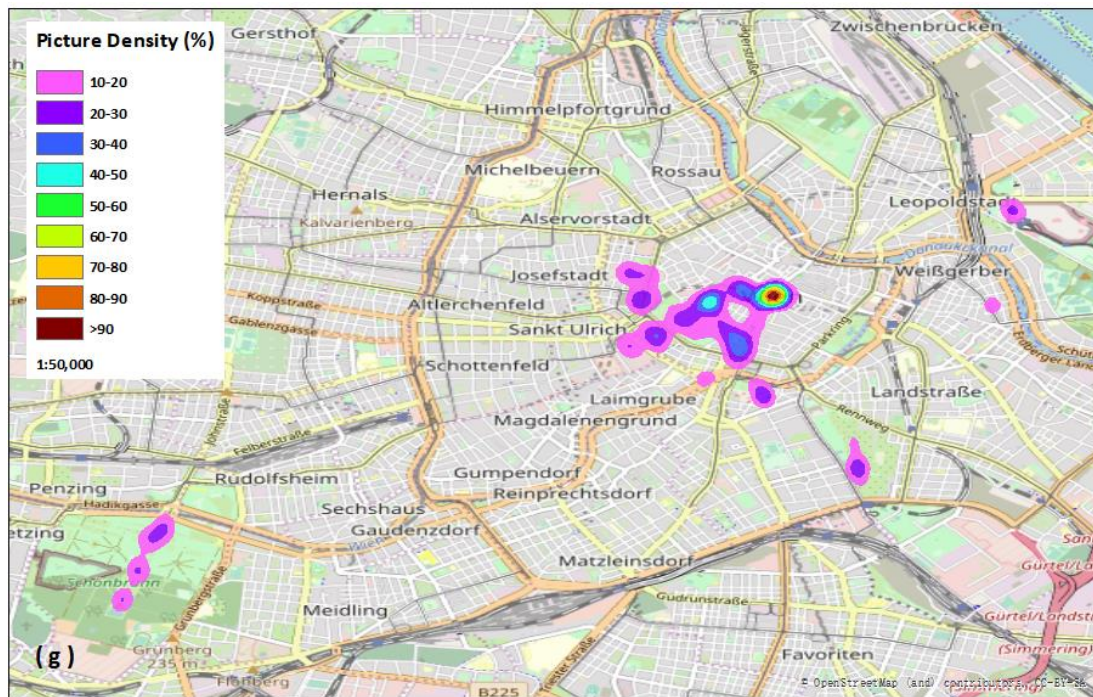


Figure 15g Footprint of the group of tourists from Italy

4.2 Modelled City Center

In this study, the city center is depicted by the smooth surface with diverse densities of city center related pictures. Figure 10 shows the resulted KDE result of pictures from locals and all tourists. The results are represented by the means of a density map. Each color represents a certain interval of the density of pictures. The higher the density value is, the more

commonly the area is considered as the city center by Flickr users. In addition, the polygons which represent thematic POIs are overlapped with the density map to be referenced. For locals, the area with a higher density of city center related pictures is more compacted. As we can see, there are only two hotspots with density higher than 60 % for local citizens; while for tourists, there are three hotspots with density higher than 60%. All of them are within the district Innere Stadt. As for the two hotspots of local citizens, the maximum picture density intervals of these two areas are 60%-70% around the intersection of Graben Street and Kohlmarkt Street, and greater than 90% close to Stephansdom (St. Stephen's Cathedral) (see Figure 11a); while for the tourists, the maximum density bands of each hotspot are 80%-90% at Michaelerplatz, 80%-90% close to Peterskirche (St. Peter's Church) and greater than 90% at St. Stephansdom (see Figure 11b). It shows that locals have a more accordant idea of defining the city center than tourists. They mostly agree that the city center is around Stephansdom. Conversely, tourists appear to be less certain about what is city center area. Although they also share the opinion with locals that the city center is mainly around Stephansdom, there are two other secondary nuclei considered to be part of the city center. Comparing the modeled city center with the extracted footprints of locals and all tourists, we can see that the area around Stephansdom where both locals and tourists mainly believe to be the city center is also the area with the highest picture density on both footprints. However, for locals, their secondary nuclei of the city center (the area around the intersection of Graben Street and Kohlmarkt Street) is not the area with second-highest picture density in their footprints. The picture density of that area is around 60%-70%. For example, three other spots around Museumsplatz, Rathausplatz and Wiener Staatsoper (Vienna State Opera) have greater than 70% density of pictures. To the opposite, for all tourists, the other two secondary nuclei of the city center are also the areas with a higher density of pictures in the footprint of all tourists. The result shows that the range of locals' activities is not restricted to the city center while tourists tend to consider the city center as their main area of activity.

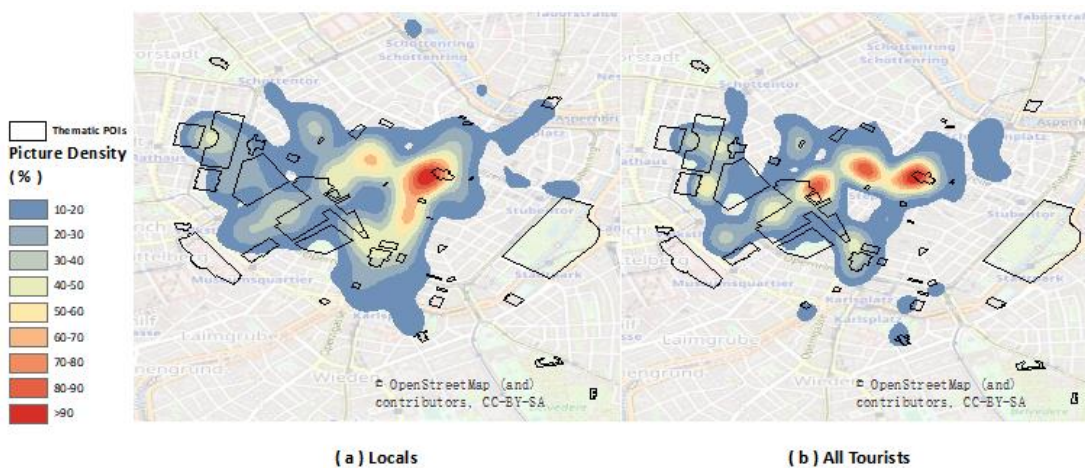


Figure 16 Modeled city center

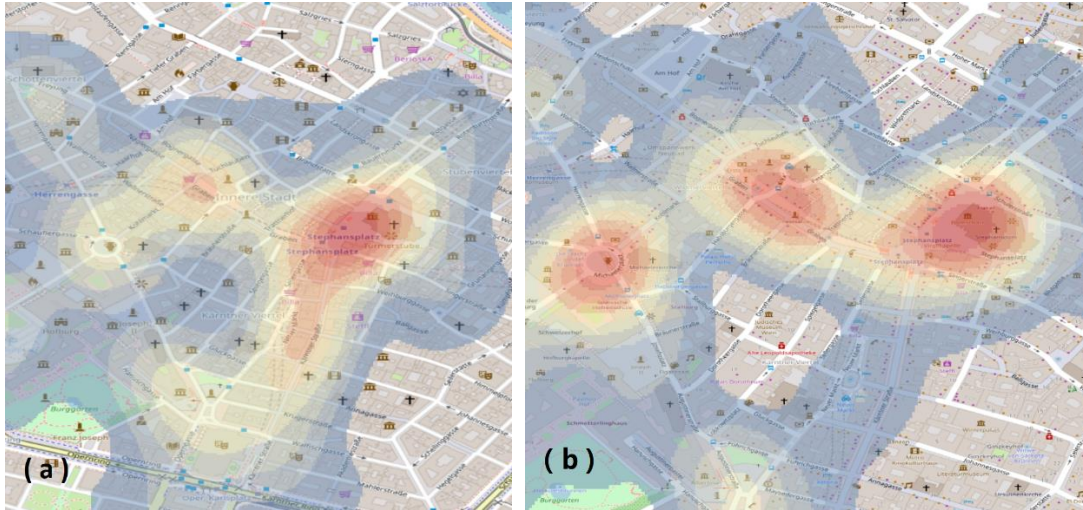


Figure 17 Zoomed city center hotspots - base map by OpenStreetMap: (a) Locals (b) All tourists

4.3 Tourist Profile

4.3.1 Thematic interests

A wind rose is used to visually represent the distinctive thematic interests of each user group during the whole period of time. Figure 12 provides a general overview of all different user groups which includes the group of local citizens, all tourists, domestic tourists, tourists from Germany, UK, US, and Italy. The axis represents the ratio of targeted POIs in each category, the higher the ratio is, the more interested the corresponding user group is to this certain category.

As we can see, none of the user groups show much interest in shopping areas as well as contemporary sights and architecture. However, there is still 20% of the shopping areas and contemporary sights are targeted by local citizens. The user group of locals has a leading position in each dimension followed by domestic tourists. And both groups have around 80% targeted POIs for religious sights as well as operas and theaters. But locals show a particular interest in museums comparing with other groups. 86% of museums are targeted by locals while only 57% are targeted by tourists from German and domestic tourists. As for the four other user groups, tourists from Germany has higher interests among historical sights and architectures, museums, nature and parks; Tourists from Germany and UK are both more attracted to operas and theaters compared with other international tourists; Tourists from Italy also has greater interest in religious and architectures; Moreover, tourists from the US shows relatively higher interest in dining and drinking places. The orange line shows the overall thematic interests of all tourists, which shows that tourists are generally more interested in museums and less interested in places for shopping, dining and drinking as well as contemporary sights and architectures.

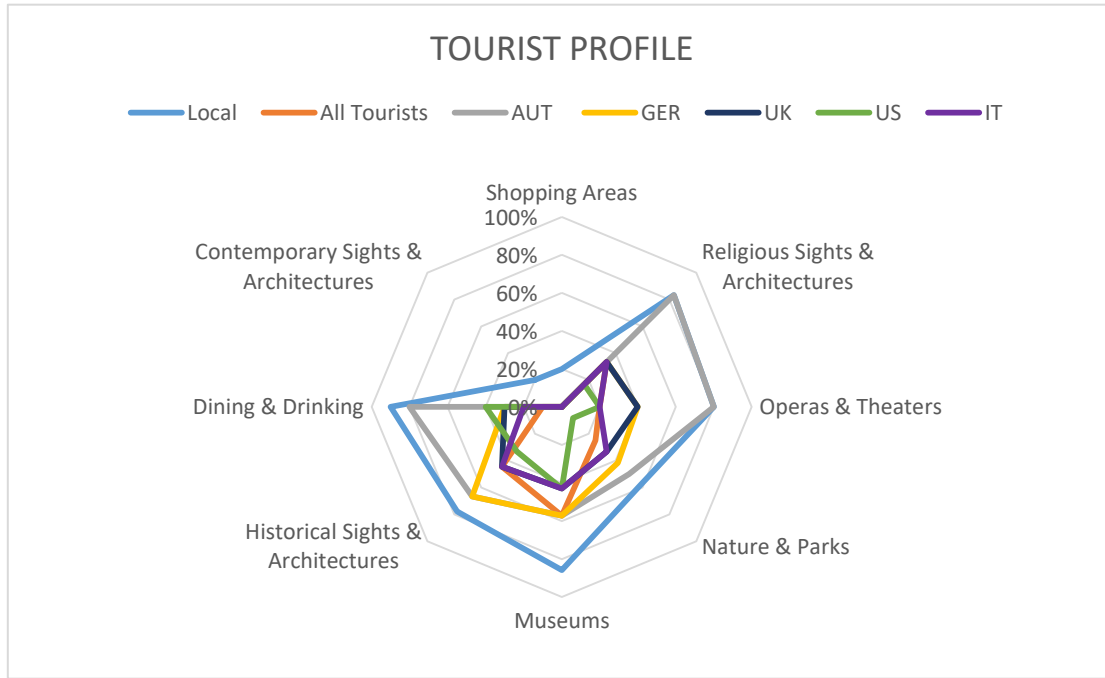


Figure 18 Tourists Profile - Overview

4.3.2 Temporal trend

The seasonal changes of each user group's footprints can be revealed by the heat map below (Appendix Figure 14). The feature value which represented the attractiveness of each POIs are represented as colors. The darker the color is, the more attractive a POI is considered. However, it is hard to conclude a general trend shared by all user groups. But some minor patterns can be noticed.

As we can see from the subplot for the user group of all tourists, there is not much change over seasons for all POIs. However, Wiener Staatsoper (Vienna State Opera), Albertina and Heldenplatz (Heroes' Square) are attracting more pictures in spring and autumn. While for locals, people tend to be more into parks and museums in autumn than any other season. For domestic tourists, POIs regarding the category "Nature & Parks" are less popular in summer and winter. In addition, for tourists from the US, they seem to be more interested in historical sights and architectures (Heldenplatz and Spanish riding school for example) in spring and autumn. And the most popular POI -- Stephansdom (St. Stephan's Cathedral), the feature value barely changes among all user groups in all the seasons.

The number of pictures in each season throughout time for each user group is displayed in Figure 13. In the line graph, pictures of every four seasons in every four years are summed up. Some trends can be found in this graph. As it is shown in the graph, the number of pictures of each season are small but remain stable over seasons from March 2002 to February 2006. The remarkable growth of the number of pictures appears in the spring of 2006 for all user groups. Since then, for the group of all tourists, summer is the season which most uploaded pictures. The number of pictures uploaded by all tourists reaches a peak in the summer of the

year 2010 to 2014. Also, for domestic tourists and locals, the number of pictures in the spring, summer, and autumn are relatively stable since the year 2006 to 2014. Also, for the group of tourists from Germany, the US, and Italy, we can notice that a lot fewer pictures are uploaded in autumn than in summer.

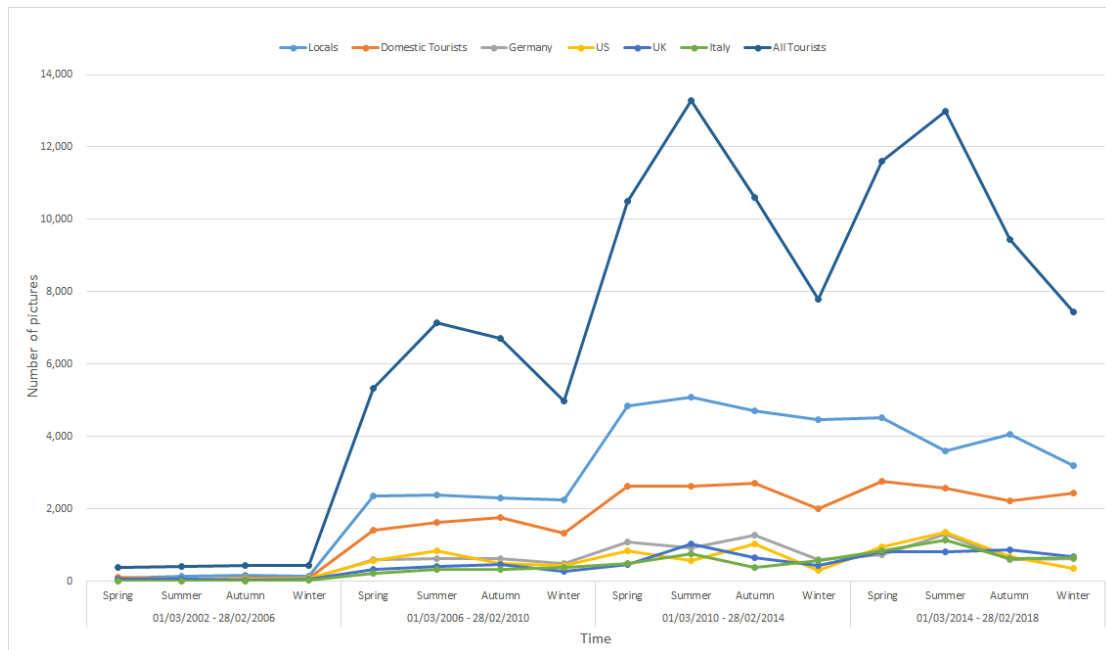


Figure 19 The number of pictures in each season

5. Discussions

As the research demonstrates, Flickr data as one type of VGI has advantages for the study of urban traces. Different from the traditional data obtained by census or statistical bureaus, the VGI obtained from social media platform provides valuable information with a remarkably finer-grained spatial and temporal resolution at low cost (Jiang, Ma, Yin, & Sandberg, 2016). Comparing with other techniques and sources of evidence, it offers a better and cheaper resource to answer this type of question which involves human perception (Elwood, Goodchild, & Sui, 2012). The motive of people contributing to VGI helps us getting a deeper insight into tourists' subjective ideas which is optimal for questions like the fuzzy concept of the city center. Despite all the advantages of using VGI, there are inevitable disadvantages of it. Deficiencies lie in several aspects.

First, as it has been widely discussed, the quality of VGI is not assured (Goodchild & Li, 2012). Unlike those traditional professional information sources which uphold high-quality standards, the credibility of VGI is to be questioned. The contributors of these data are mostly general public instead of experts or scientists. For the geotagged Flickr data, there are possibilities of false location tagging. As mentioned in the data review, pictures without auto-geotagging require manual operations. The localness of the VGI contributor is important (Johnson, Sengupta, Schöning, & Hecht, 2016). Since most pictures are uploaded by tourists

who are unfamiliar with the city, it is much likely for them to assign false locations to their pictures. By comparison, the location information uploaded by local citizens who are immersed in the urban environment seems to be more reliable. Nevertheless, even if the users are well aware of the environment, there are still chances for them to commit errors. As Marlow mentioned, “user incentives and motivations may influence the resultant tags in a tagging system” (Marlow, Naaman, Boyd, & Davis, 2006). When people uploading pictures on social media platforms like Flickr, what they focus on is the image instead of the precise location of taking pictures. Social media users might just tag an approximate location according to their vague memories. For them, the accuracy of the geotag is of less importance.

The characteristic of data is largely dependent on the generators of it. More specifically, the representativeness of VGI is limited due to the uncertainties in the demographics of social media users (Kovacs-Györi et al., 2018). The issues of the “digital divide” (Wiersma, 2010) lead to the uniformity of the Flickr user group. People who are with less digital literacy are under-represented in the Flickr dataset. Groups with very young ages and senior ages are less likely to be included in the dataset. Additionally, since the data is only acquired from Flickr, only Flickr users are studied in this research. It is obvious that not all tourists are active with social media. And even if they are social media users there are multiple other choices of social media platforms. According to the tourism statistics of Vienna (“Vienna Tourist Board: Arrivals & bednights 2018”, 2019), China is one of the main tourism markets (ranked the 7th place) for Vienna. However, depending on the ranking of the number of pictures on Flickr, China is far behind other countries. Therefore, the preference for using social media impacts greatly on the representativeness of VGI as well.

Another issue that might introduce bias is the classification of locals and tourists. It is simple and direct to classify those Flickr users who have valid information reveal their localness. But for users who have ambiguous origin information, the localness is determined by the classification based on the temporal feature of their Flickr uploads. The method has uncertainties. Because the behavior of uploading Vienna-related posts does not 100 percent reveals users’ actual activities in Vienna. For instance, some users might tend to upload pictures which are taken on different dates at one time after their visit.

Furthermore, whether and how a POI will be photographed is dependent on the type of the place and the type of activities people conduct at this spot as well. Some types of POIs will be underestimated due to the lack of uploaded pictures. But it does not necessarily mean that in reality, people do not visit those places as much as others. For example, people might not tend to take pictures when they are having coffee in the Café or shopping in the shopping mall. Also, it is sometimes not possible to get access to the Internet at some places (spots located in the mountain area for example). So people cannot upload the related pictures in time. The delay in uploading might lead to a higher chance of false geotagging. While for some POIs, the location of the uploaded pictures is mostly not exactly where the POI is located. POIs like museums, architectures, operas and theaters are always photographed from the outside of it. So the concentrated spot of related pictures is dislocated. Although this is not a part of this study, new insights could be provided regarding the interpretation of footprints.

6. Conclusion

In this study, an approach was designed to differentiate the urban traces left by tourists from diverse origins and local citizens based on the volunteered geographic information. The approach is tested on the Flickr data in the city of Vienna. It has been proved that this approach can achieve the goal of differentiating urban traces. Furthermore, it provides a deeper insight into tourists as well as local citizens' visiting behavior and the concept of places in the urban environment.

The research questions addressed in chapter 1.2.2 can be answered as follows:

Sub-objective a: To map the urban traces of tourists and local citizens from social media presented by their distinctive footprints

Research question: Are there differences in footprints between tourists from different origins and local citizens? Which are those differences?

As it is shown in the result chapter, distinctive footprints are generated for each user group. Despite the similarity in patterns, there are certainly differences among these footprints. The footprint of the local user group is more dispersed and covers a larger area. There are more hotspots (areas with relatively higher picture density) in the local footprints. As for the footprints of domestic tourists, it shares some features from both the footprints of locals and international tourists. For example, comparing with other international tourists, in both locals and domestic tourists' footprints, Schönbrunn zoo occupies a higher concentration of pictures while Prater occupies relatively lower concentration. And for the other study user groups (tourists from Germany, US, UK and Italy), each of their footprints have their own features at different spots. For instance, there is a higher concentration of pictures at Prater in the footprints of tourists from Germany and the UK; whereas there is a lower picture density at Belvedere and Prater in the footprints of tourists from the US.

Sub-objective b: To model the city center according to the semantics extracted from VGI of tourists and local citizens

Research question 1: How differently do tourists and local citizens perceive the city center?

As the modeled city center depicts, tourists and local citizens do have a rough agreement on the location of the fuzzy concept – city center. They both consider that area around Stephansdom is the city center and all the candidate locations are located in the district Innere Stadt. But as we can see from the result, local citizens obviously have a clearer idea of the range of the city center. Conversely, tourists seem to be more ambiguous about it. All three nuclei locations share almost the same density of city center-related pictures.

Research question 2: Is there a relation between the footprints and perceived city center among tourists and local citizens? Is this relation clearer among certain user groups?

The answer is positive. The area of Stephansdom, where both groups believe the city center is, shows a picture density greater than 90% in the footprints of both groups. However, the relation is clearer among the tourist user group. The perceived city center of tourists coincides with the area of higher picture densities in their footprint. It can be inferred that the main area of tourist's urban traces is restricted to where they consider being the city center. But locals have a wider range of visiting the city, so in their footprints, the area of higher picture density is extended.

Sub-objective c: To create a tourists profile categorized by the origin countries of tourists as well as the local citizens in respect of diverse thematic POIs

Research question 1: Can we identify a unique tourist profile regarding different thematic POIs for different user groups?

Through the last data analyzing phase, the goal of generating a tourist profile has been achieved by comparing thematic POIs with extracted AOIs for each user group. Each group shares its own emphasis of interest on different thematic categories.

Research question 2: Are there correlations between the targeted thematic POIs in the diverse footprints and specific origin countries? Is there a seasonal trend among them?

As we can see from the wind rose graph, the targeted POIs in each footprint are correlated with the origin of user groups. Different thematic POIs are targeted for each group. For example, local citizens show a relatively greater interest in all categories. They own the largest ratio of targeted POIs of all eight categories. Domestic tourists rank second. However, in the footprint of domestic tourists, POIs of categories except for religious sights and architecture as well as opera and theaters are less targeted. Tourists from Germany shows a higher interest of visiting historical sights and architectures than other international tourist groups; dining and drinking POIs seem to be more popular with tourists from the US; tourist from Italy shows least interests in dining and drinking spots; tourists from Germany and UK are more interested in operas and theaters. However, the least amount of contemporary sights and architectures, as well as shopping areas, are targeted in all the footprints of diverse tourist groups.

As for the seasonal trend, no obvious general trends can be concluded from the heat map. However, it is obvious that the feature value of the most popular POI – Stephansdom remains high in all groups throughout time. And also, by analyzing the number of uploaded pictures, we can see that for the group of all tourists, the largest number of pictures are always uploaded in summer.

In conclusion, the adapted approach is able to differentiate urban traces of tourists and locals methodologically. Distinctive patterns of different social media user groups can be revealed by the digital footprints of their uploaded pictures. Comparing with geographic data obtained by traditional approaches and data sources, VGI does not only offer a wider spatial extent of data with a finer spatial and temporal resolution, but it also provides better solutions to study human perception related problems like the fuzzy concept of places. Humans are utilized as sensors, which directly indicates how humans perceive the outer environment. However, traditional surveys and data collection cannot be replaced by VGI. The deficiencies of VGI like the lack of reliability of data and the under-representation of user groups introduce bias to the study results inevitably. The tourism official data is undoubtedly more accurate about information like the population of tourists from different origins and temporal features of visiting certain cities. Also, surveys about the intention of tourists' visiting behavior can directly reveal their diverse interests regarding different thematic types of sights. The main obstacles of social media data analysis are still the traditional disadvantages of VGI. Even though, this approach can provide some insights about urban tourism through grasping general patterns as the foundation of further in-depth analysis and field research for relevant experts (Kovacs-Györi et al., 2018).

REFERENCE

- Adrianne Jeffries. (2013). The man behind Flickr on making the service 'awesome again'. Retrieved from <https://www.theverge.com/2013/3/20/4121574/flickr-chief-markus-spiering-talks-photos-and-marissa-mayer>
- Appreciating Vienna. Retrieved 21 August 2019, from <https://www.wien.info/en/shopping-wining-dining>
- Bertini, E., Di Girolamo, A., & Santucci, G. (2007). *See What You Know: Analyzing Data Distribution to Improve Density Map Visualization*. Paper presented at the EuroVis.
- Bevölkerung zu Jahres-/Quartalsanfang. (2019). Retrieved 2019, August 14, from http://www.statistik.at/web_de/statistiken/menschen_und_gesellschaft/bevoelkerung/bevoelkerungsstand_und_veraenderung/bevoelkerung_zu_jahres-_quartalsanfang/023582.html
- Boeing, G. (2018). Clustering to Reduce Spatial Data Set Size. Retrieved from <https://dx.doi.org/10.31235/osf.io/nzhdc>. doi:10.31235/osf.io/nzhdc
- Buhalis, D., & Amaranggana, A. (2015). Smart tourism destinations enhancing tourism experience through personalisation of services. In *Information and communication technologies in tourism 2015* (pp. 377-389): Springer.
- Coleman, D., Georgiadou, Y., & Labonte, J. (2009). Volunteered geographic information: The nature and motivation of producers. *IJSDIR*, 4(1), 332-358.
- Donaire, J. A., Camprubí, R., & Galí, N. (2014). Tourist clusters from Flickr travel photography. *Tourism management perspectives*, 11, 26-33.
- Dotan, A., & Zaphiris, P. (2010). A cross-cultural analysis of Flickr users from Peru, Israel, Iran, Taiwan and the United Kingdom. *International Journal of Web Based Communities*, 6(3), 284-302.
- Elwood, S., Goodchild, M. F., & Sui, D. Z. (2012). Researching volunteered geographic information: Spatial data, geographic research, and new social practice. *Annals of the association of American geographers*, 102(3), 571-590.
- Feick, R., & Roche, S. (2013). Understanding the Value of VGI. In *Crowdsourcing geographic knowledge* (pp. 15-29): Springer.
- Flanagin, A. J., & Metzger, M. J. (2008). The credibility of volunteered geographic information. *GeoJournal*, 72(3-4), 137-148. Retrieved from <https://dx.doi.org/10.1007/s10708-008-9188-y>. doi:10.1007/s10708-008-9188-y
- Fonte, C., Bastin, L., Foody, G., Kellenberger, T., Kerle, N., Mooney, P., . . . See, L. (2015). VGI quality control. *ISPRS Geospatial week 2015*, 317-324.
- Gao, S., Janowicz, K., & Couclelis, H. (2017). Extracting urban functional regions from points of interest and human activities on location-based social networks. *Transactions in GIS*, 21(3), 446-467. doi:10.1111/tgis.12289
- Gao, S., Janowicz, K., Montello, D. R., Hu, Y., Yang, J.-A., McKenzie, G., . . . Yan, B. (2017). A data-synthesis-driven method for detecting and extracting vague cognitive regions. *International Journal of Geographical Information Science*, 31(6), 1245-1271.
- García-Palomares, J. C., Gutiérrez, J., & Mínguez, C. (2015). Identification of tourist hot spots based on social networks: A comparative analysis of European metropolises using photo-sharing services and GIS. *Applied Geography*, 63, 408-417.
- Girardin, F., Calabrese, F., Dal Fiore, F., Ratti, C., & Blat, J. (2008). Digital footprinting: Uncovering tourists with user-generated content. *IEEE Pervasive computing*, 7(4), 36-43.

- Girardin, F., Vaccari, A., Gerber, A., Biderman, A., & Ratti, C. (2009). Quantifying urban attractiveness from the distribution and density of digital footprints.
- Goodchild, M. F. (2007). Citizens as sensors: the world of volunteered geography. *GeoJournal*, 69(4), 211-221.
- Goodchild, M. F., & Li, L. (2012). Assuring the quality of volunteered geographic information. *Spatial statistics*, 1, 110-120.
- Grothe, C., & Schaab, J. (2009). Automated Footprint Generation from Geotags with Kernel Density Estimation and Support Vector Machines. *Spatial Cognition & Computation*, 9(3), 195-211.
doi:10.1080/13875860903118307
- Haklay, M. (2010). How good is volunteered geographical information? A comparative study of OpenStreetMap and Ordnance Survey datasets. *Environment and planning B: Planning and design*, 37(4), 682-703.
- Haklay, M. (2013). Citizen Science and Volunteered Geographic Information: Overview and Typology of Participation. In (pp. 105-122): Springer Netherlands.
- Hasan, S., Zhan, X., & Ukkusuri, S. V. (2013). *Understanding urban human activity and mobility patterns using large-scale location-based data from online social media*. Paper presented at the Proceedings of the 2nd ACM SIGKDD International Workshop on Urban Computing - UrbComp '13.
- Hatz, G. (2008). Vienna. *Cities*, 25(5), 310-322.
- Hu, Y., Gao, S., Janowicz, K., Yu, B., Li, W., & Prasad, S. (2015). Extracting and understanding urban areas of interest using geotagged photos. *Computers, Environment and Urban Systems*, 54, 240-254.
doi:10.1016/j.compenvurbsys.2015.09.001
- Jansen-Verbeke, M. (1992). Urban recreation and tourism: physical planning issues. *Tourism Recreation Research*, 17(2), 33-45.
- Jiang, B. (2013). Volunteered geographic information and computational geography: New perspectives. In *Crowdsourcing geographic knowledge* (pp. 125-138): Springer.
- Jiang, B., Ma, D., Yin, J., & Sandberg, M. (2016). Spatial distribution of city tweets and their densities. *Geographical Analysis*, 48(3), 337-351.
- Johnson, I. L., Sengupta, S., Schöning, J., & Hecht, B. (2016). *The geography and importance of localness in geotagged social media*. Paper presented at the Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems.
- Kádár, B. (2013). Differences in the spatial patterns of urban tourism in Vienna and Prague. *Urbani izziv*, 24(2), 96-111.
- Kernel Density. Retrieved 2019, August 21 from <http://desktop.arcgis.com/en/arcmap/10.3/tools/spatial-analyst-toolbox/kernel-density.htm>
- Kovacs-Györi, A., Ristea, A., Kolcsar, R., Resch, B., Crivellari, A., & Blaschke, T. (2018). Beyond Spatial Proximity—Classifying Parks and Their Visitors in London Based on Spatiotemporal and Sentiment Analysis of Twitter Data. *ISPRS International Journal of Geo-Information*, 7(9). doi:10.3390/ijgi7090378
- Kriegel, H. P., Kröger, P., Sander, J., & Zimek, A. (2011). Density - based clustering. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1(3), 231-240. doi:10.1002/widm.30
- Mark, J. (1939). The Law of the Primate City. *Geographical Review*, 29(2), 226-232.
- Marlow, C., Naaman, M., Boyd, D., & Davis, M. (2006). *HT06, tagging paper, taxonomy, Flickr, academic article, to read*. Paper presented at the Proceedings of the seventeenth conference on Hypertext and hypermedia.
- McKenzie, G., & Adams, B. (2017). *Juxtaposing thematic regions derived from spatial and platial user-generated content*.

- Montello, D. R., Friedman, A., & Phillips, D. W. (2014). Vague cognitive regions in geography and geographic information science. *International Journal of Geographical Information Science*, 28(9), 1802-1820.
- Montello, D. R., Goodchild, M. F., Gottsegen, J., & Fohl, P. (2017). Where's downtown?: Behavioral methods for determining referents of vague spatial queries. In *Spatial Vagueness, Uncertainty, Granularity* (pp. 185-204): Psychology Press.
- O'sullivan, D., & Unwin, D. (2014). *Geographic information analysis*: John Wiley & Sons.
- Popescu, A., & Grefenstette, G. (2009). *Deducing trip related information from flickr*. Paper presented at the Proceedings of the 18th international conference on World wide web.
- 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.
- seaborn.heatmap. Retrieved 2019, August 21, from <https://seaborn.pydata.org/generated/seaborn.heatmap.html>
- Sightseeing in Vienna. Retrieved 2019, August 21, from <https://www.wien.info/en/sightseeing>
- Stefanidis, A., Crooks, A., & Radzikowski, J. (2013). Harvesting ambient geospatial information from social media feeds. *GeoJournal*, 78(2), 319-338.
- Sui, D., Elwood, S., & Goodchild, M. (2012). *Crowdsourcing geographic knowledge: volunteered geographic information (VGI) in theory and practice*: Springer Science & Business Media.
- Sun, Y., Fan, H., Helbich, M., & Zipf, A. (2013). Analyzing Human Activities Through Volunteered Geographic Information: Using Flickr to Analyze Spatial and Temporal Pattern of Tourist Accommodation. In J. M. Krisp (Ed.), *Progress in Location-Based Services* (pp. 57-69). Berlin, Heidelberg: Springer Berlin Heidelberg.
- Verstockt, S., Milleville, K., Ali, D., Porras-Bernardez, F., Gartner, G., & Van de Weghe, N. (2019). *EURECA - European Region Enrichment in City Archives and collections*. Paper presented at the DACH'19, the 14th ICA Conference
- Vienna ranked as most liveable city in the world. (2018, August 14). Retrieved 2019, August 14, from <https://www.bbc.com/news/business-45174600>
- Vienna Tourist Board: Arrivals & bednights 2018. Retrieved from <https://b2b.wien.info/de/statistik/daten/naechtigungen-2018>
- Wiersma, Y. (2010). Birding 2.0: Citizen Science and Effective Monitoring in the Web 2.0 World* Ornithologie 2.0: la science citoyenne et les programmes de suivi à l'ère d'internet 2.0. *Avian Conservation and Ecology*, 5(2), 1-9.
- Yin, Z., Cao, L., Han, J., Luo, J., & Huang, T. (2011). *Diversified trajectory pattern ranking in geo-tagged social media*. Paper presented at the Proceedings of the 2011 SIAM International Conference on Data Mining.
- Zhang, Y. (2019). An approach to localness assessment of social media users.

APPENDIX

Table 5 List of thematic POIs

Categories	Name of the POIs	Categories	Name of POIs
Shopping Area (Total: 10)	Lugner City GmbH Wien Mitte The Mall Ringstraßen-Galerien Galleria Columbus Center Donau Zentrum Shopping Center Nord Zentrum Simmering Einkaufszentrum Hernals Anhof Center	Museums (Total: 7)	Mozarthaus Vienna Kunsthistorisches Museum MuseumsQuartier Vienna Secession House of Music Albertina Upper Belvedere
Religious Sights & Architectures (Total: 6)	Minoriten Church St.Stephen's Cathedral St.Charles' Church Synagogue St. Rupert's Church Votice Church	Historical Sights & Architectures (Total: 9)	Imperial Palace & Heldenplatz Parliament Vienna city hall Am Hof Square Judenplatz Schoenbrunn Palace Lower Belvedere Hunderwasser House Spanish Riding School
Natures & Parks (Total: 12)	Volksgarten Türkenschanzpark Stadtpark Kurpark Oberlaa Augarten Burggarten CityhallPark Park at Schönbrunn Palace Alpengarten Schönbrunn Zoo Prater Donaupark	Dining & Drinking (Total: 10)	Café Central Café Frauenhuber Gerstner K. & K.Hofzuckerbäckerei Café Hawelka Café Imperial Café Mozart Café Museum Café Schwarzenberg Conditorei Sluka Heuriger am Belvedere
Operas & Theaters (Total: 5)	Wiener Staatsoper Musikverein Vienna Konzerthaus Burgtheater Volkstheater	Contemporary Sights & Architectures (Total: 5)	Danube Tower DC tower Vienna's Gasometers SO/Vienna Campus WU (library)

Table 6 City center related tags in multi-languages

Language Code	Tags	English Translation	Language Code	Tags	English Translation
EN	center	-	DE	zentrum	center/centre
EN/FR	centre	-	DE	stadtzentrum	citycenter
EN	city center	-	DE	stadt zentrum	city center
EN	citycenter	-	DE	innerestadt	innercity
EN	central	-	DE	innere stadt	inner city
EN	downtown	-	DE	innenstadt	citycenter
EN	inner city	-	DE	innen stadt	city center
EN	innercity	-	DE	stadtkern	urban core
EN	CBD	-	DE	stadttinneres	city center
EN	urban core	-	DE	alterstadt	old town
EN	urbancore	-	ES/IT/PT	centro	center
EN	old town	-	FR	centre-ville	downtown
EN	oldtown	-	FR	ville-centre	city center
EN	quartier central	-	FR	en ville	downtown
HU	belváros	downtown	RU	центр	center
HU	centrum	center	AR	مركز	center
HU	középpont	center	AR	المركز	downtown
HU	városközpont	city center	SQ	qendër	center
HU	Óváros	old town	JA	ダウンタウン	downtown
ZH	市中心	city center	JA	センター	center
ZH	中心	center	JA	市の中心部	city center
ZH	老镇	old town	ZH	老城	old town

Figure 20 Heatmap

