



**POLITECNICO**  
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE  
E DELL'INFORMAZIONE

# Location Inference through Social Media and Social Relationships

TESI DI LAUREA MAGISTRALE IN  
COMPUTER SCIENCE AND ENGINEERING - INGEGNERIA IN-  
FORMATICA

Author: **Matteo Rizzi**

Student ID: 10535464

Advisor: Prof. Stefano Longari

Co-advisors:

Academic Year: 2022-23



## Abstract

As of today, social media occupy a big role in our lives, most people have one or more social media accounts where they post regularly documenting their lives. This has brought to our society various improvements in how we connect with people and overall in our quality of life, but what are the negative consequences? Continuously sharing information about our life on social media can reveal many sensitive information about ourselves. An important one is our location. Location inference is the practice of discovering someone's location without them disclosing it directly. For these reasons, data from social media can be used to infer someone's location, without their consent. There are many ways to achieve this: by using mentioned location in posts, local words and, more importantly, activity of a user's friends on social media. Thanks to those information, it is possible to infer a user's location even through social media activity that is not coming directly from them. The intricate interplay between online social connections and physical geography has catalyzed the emergence of innovative strategies for predicting user locations. Many works till today have studied this phenomenon and proposed their methods to infer someone's location through their friends as well as other sources of information. Because of this, a work that organizes them in order to give a complete overview about this subject is needed. That is the reason of this survey: to examine a collection of methodologies, algorithms, and models that exploit interconnected relationships to infer someone's location leveraging social media data. While the scope of this survey is to give a complete overview about the state of the art regarding location inference through friendships, in order to give a more complete and broader description of the subject this survey also covers the various aspects related to location data and location privacy. To achieve this, it also analyzes the topic of location privacy protection mechanisms, location inference in general, and friendship inference using location data from social media.

**Keywords:** location inference, location data in social media, location privacy, user's friendships



## Abstract in lingua italiana

Al giorno d'oggi, i social media svolgono un grande ruolo nelle nostre vite. La maggior parte delle persone possiede uno o più account sui social media, dove condivide regolarmente fatti riguardanti la propria vita. Questo ha introdotto nella nostra società numerosi miglioramenti per la qualità della vita e il modo in cui restiamo in contatto con altre persone. Ma quali sono le conseguenze negative? Condividere continuamente informazioni relative alla vita privata può rivelare molte informazioni sensibili. Una di esse è la posizione. Inferire la posizione di qualcuno consiste nello scoprire la sua posizione senza che questi la riveli intenzionalmente. Siccome i social media ricoprono un ruolo significativo nelle nostre vite, come detto, e siccome rilasciamo una quantità rilevante di informazioni ogni giorno su queste piattaforme, esse sono un'ottima fonte di informazioni per inferire la posizione di qualcuno. Ci sono molti modi di scoprire la posizione di un utente: sfruttando i posti menzionati, le parole tipiche di alcune regioni geografiche e, soprattutto, sfruttando le attività degli amici sui medesimi social. Grazie a queste informazioni, è appunto possibile inferire la posizione di un utente grazie ad attività sui social non provenienti direttamente da lui. L'intricato intreccio fra relazioni sociali sui social media e le posizioni geografiche ha contribuito allo sviluppo di strategie innovative per inferire le posizioni degli utenti. Molti lavori hanno studiato questo fenomeno e proposto i loro metodi per inferire la posizione di qualcuno tramite le sue amicizie e altre informazioni tratte dai social media. Per questo motivo serve una survey che li organizzi al fine di fornire una panoramica completa riguardo questo campo. E' questa la motivazione alla base di questo lavoro: esaminare una serie di metodi, algoritmi, e modelli che sfruttano le interconnessioni tra le persone per scoprire la posizione di qualcuno attraverso i social media. Nonostante lo scopo principale di questa tesi sia fornire una panoramica dello stato dell'arte sull'inferire la posizione di un utente attraverso le sue amicizie sui social media, per fornire una descrizione completa e ampia del problema questa tesi tratta anche i vari aspetti della localizzazione e della privacy relativa alla posizione. Per ottenere questo risultato, questo lavoro analizza anche i vari meccanismi per proteggere la privacy della posizione, analizza il problema della localizzazione in generale, e anche il problema dell'inferenza delle amicizie nella vita reale tramite i dati della posizione ricavabili dai

social media.

**Parole chiave:** localizzazione, dati di posizione nei social media, privacy della posizione, amicizie fra utenti

# Contents

<b>Abstract</b>	<b>i</b>
<b>Abstract in lingua italiana</b>	<b>iii</b>
<b>Contents</b>	<b>v</b>
<b>Introduction</b>	<b>1</b>
<b>1 Related Works</b>	<b>5</b>
<b>2 Location Privacy Protection</b>	<b>9</b>
2.1 Location Privacy Protection Mechanisms . . . . .	9
2.2 Privacy evaluation metrics . . . . .	18
2.3 Effectiveness of LPPMs . . . . .	21
<b>3 Location Inference</b>	<b>23</b>
3.1 Home Location Inference . . . . .	24
3.2 Activity Location Inference . . . . .	25
3.3 Next Location Prediction . . . . .	28
3.4 Location Inference through Social Network . . . . .	32
3.4.1 Location Inference through Friendships . . . . .	32
3.4.2 Location Inference through User Similarity . . . . .	41
<b>4 Friendship Inference through Location Data</b>	<b>49</b>
4.1 Friendship Inference through Check-Ins . . . . .	49
4.2 Friendship Inference through Trajectories . . . . .	51
4.3 Friendship Inference through Multiple Sources Including Location Data . .	53
<b>5 Discussion and Open Research Questions</b>	<b>55</b>
<b>6 Conclusions</b>	<b>59</b>

<b>Bibliography</b>	<b>61</b>
<b>List of Figures</b>	<b>75</b>
<b>List of Tables</b>	<b>77</b>



# Introduction

Discussing how location privacy can be inferred from social media data is an important subject due to the increasing integration of location-based services and the pervasive nature of digital technologies in our lives. Location data are very important due to many reasons: locating data is highly sensitive and can reveal intimate details about individuals' lives, such as their daily routines, preferences, and social relationships. Location data can be used to profile users basing on their habits, preferences, and demographics. Moreover social media data from people who we interact with can affect our privacy. Many factors today make talking about this subject more and more relevant: the advancement of technology increases the potential for more accurate and granular location data inference; thanks to the widespread use of smartphones and the integration of location-sharing features in social media platforms, a lot of the population is now disclosing its whereabouts. And even when users don't disclose their position directly, it is still possible to infer it by analyzing the other data that users post on social platforms. In fact it is possible to infer someone's location by using mentioned location in posts, local words and activity of a user's friends on social media. This way it is possible to infer users' location even through social media activity that is not coming directly from them. This correlation between online social relationships and physical locations has caused the emergence of new and innovative ways to infer someone location. As a matter of fact, numerous studies to date have explored the concept of inferring individuals' location by leveraging their social connections and various sources of data. Given the diversity of these approaches, a comprehensive survey is necessary to synthesize and present these methodologies. This survey aims to fulfill this need by systematically examining a range of techniques, algorithms, and models that exploit interconnected relationships to deduce a person's location from the data available on social media platforms. While the primary goal of this survey is to offer a comprehensive overview of the latest advancements in location inference through friendships, its scope extends beyond that. It encompasses a broader understanding of location-related topics, including the various aspects of location data and location privacy. In fulfilling this comprehensive approach, the survey also delves into subjects such as safeguarding location privacy, broader location inference methodologies, and the utilization

of location data in deducing social connections. By embarking on this comprehensive exploration, the survey not only provides insights into the state of the art in location inference through social connections but also offers a well-rounded understanding of the landscape of location data and its privacy implications.

Given the abundance of works in this field, this thesis focuses on giving an overview by collecting and categorizing them. This to give a base for future works and researches that will bring innovation into this field. This thesis will cover techniques that use data from centralized social media, it will not include distributed tracking applications such as the ones used for the pandemic of COVID-19. It will cover data from social media such as Twitter, Facebook etc., and will also cover data from more location-centric social media and mobile applications such as Foursquare. The methods that this thesis covers exploit various sources such as check-ins, location trajectories, text of post in social media, following relationships on social media, mentions of other users and other interactions, etc.

The main aims of this work are the following:

- Giving an overview of the defense mechanisms for preserving location privacy.
- Giving an exhaustive overview of how social media's data and data from other services can be used to infer location of users, especially when adding information regarding the social network of a user.
- Covering the subject of inferring social relations using location data from social media and other services.

To achieve to goals mentioned before, this work is structured as follows: In Chapter 1, this work will cover works that focused on the same or similar subjects covered in this thesis. This by describing their focus, achievements, and differences from this survey. The following Chapter 2 will start by covering the topic about location privacy preserving mechanisms, highlighting their characteristics, how they work, how to evaluate them, and how effective they are. This section is particularly important for the aim of this thesis because it is necessary to cover the protection mechanisms and privacy models that are protecting our sensitive data today. This section comes before the section about location inference. The reason behind this choice is to give the reader a general understanding of the protection mechanisms the location inference attacks covered in the following chapters will have to deal with. So, in Chapter 3, the topic of location inference will be covered, dividing the examined works in three main categories based on what type of location is being inferred: home location, activity location and next location. Then, Chapter 3, will

cover the topic of inferring location through users' social network (users relationships such as friends, and people surrounding them), by differentiating methods that uses data from friends from methods that exploit similarity between users' behaviors. In Chapter 4, the subject of inferring friendships using location data will be briefly covered. This topic is covered after the previous one because it covers the opposite side of the problem covered in Chapter 3 (location inference through friends' social media) and it is interesting to see the correlations between them both. Then, in Chapter 5, this work will present its observations and open research questions. Finally, Chapter 6 will contain the conclusions.



# 1 | Related Works

Location inference and privacy in social media is a broad subject that touches many fields. To define it, location inference is a particular type of attribute inference. With attribute inference we refer to the process of deducing sensitive or personal information about individuals through the analysis of their publicly available attributes or characteristics. These attributes could include demographic information, behavioral patterns, preferences, or even seemingly innocuous data points. The goal of attribute inference is to unveil additional information beyond what is explicitly disclosed by an individual. For instance, if someone's age and educational background are publicly known, attribute inference might involve deducing their income level or political affiliations. This process can be carried out through data mining, statistical analysis, machine learning techniques, or a combination of these methods. About the specific topic of this work, location inference is a category of attribute inference, where the inferred attribute is someone's geographical location. So we can define location inference as the process of deducing or predicting individuals' physical location based on the analysis of various sources of information, such as their online activities, social interactions, and publicly available data. This process involves using clues or patterns in the data to estimate where a person is or has been at a particular time.

There are works that focus on covering the whole attribute inference subject, and by doing that they also cover privacy problem related to location inference. One of them is the work *Privacy Inference Attack Against Users in Online Social Networks: A Literature Review* [79]. In this work many papers are collected and discussed dividing them into different categories. This paper is especially interesting because it has a section discussing geographic location inference attacks. This section takes especially into consideration works that exploit social relationships and friends to better infer location of users. Another interesting section is the one regarding social relationships inference attacks, where some of the discussed works also use location to infer relationships. Another work that covers the whole privacy inference subject is the work *A survey on privacy in social media: Identification, mitigation, and applications* [9]. This research paper provides an overview of significant advancements in the realm of user privacy within social media. Specifically,

it comprehensively examines and contrasts cutting-edge algorithms pertaining to both privacy attacks and anonymization. The paper also delves into the diverse spectrum of privacy risks that stem from social media, categorizing them methodically. A particular section of this study is dedicated to exploring the intersection of social media users' location and privacy concerns.

Other works focus more on location inference only instead of attribute inference in general, such as the work *Mining location from social media: A systematic review* [102]. In this review the authors focus on differentiating every possible different source of data that can be used for location inference. For example they have a section for works that use videos, and a dedicated section for works that use links to other social media platforms. So they dedicate a lot of efforts to collect many works related to location inference and to categorize them using very specific categories. In the work *A survey on next location prediction techniques, applications, and challenges* [35] they focus on listing all the works that have been conducted regarding next location prediction, categorizing them in different categories. The work *A survey of location prediction on twitter* [121] centers on location prediction challenges within Twitter. In this segment, the authors commence by providing an insight into the Twitter platform. Approaching it from a typical user's perspective, they present a synopsis of the Twitter dataset, examining it through content, network, and contextual lenses. Subsequently, the work delves into three prevalent geolocation issues. These prediction challenges hinge on the previously mentioned information as their primary input. They also cover the evaluation metrics utilized. This paper is particularly interesting for its section about location inference with friendship-based methods. Another work *Location prediction in large-scale social networks: an in-depth benchmarking study* [3] categorize many location prediction models, including some that use friendships, and it highlights the most important and influential ones. The paper *An overview of microblog user geolocation methods* [64] presents a comprehensive analysis of users' geolocation techniques in the microblogging domain. The aim of the paper is to systematically assess, categorize, and juxtapose the efficacy of current methods for geolocating users within microblogs. The authors offer a comprehensive user geolocation framework, summarizing approaches outlined in prior literature, and succinctly outlining the merits and constraints of these methods. Their assessment is based on a performance evaluation achieved by comparing experimental outcomes reported in existing literature, all on the same datasets. Again, this work is particularly interesting for its section about network-based methods.

There are works that focus more on the exploitation of a user social network and social relationships to better infer the location. The paper *Geolocation prediction in twitter using*

*social networks: A critical analysis and review of current practice* [52] focuses on reviewing the state of the art's methods for location inference using the social network. In particular, they tested nine geolocation inference techniques. The main contributions of this thesis are: firstly, the authors underscore that the practical performance of algorithms often falls short of the accuracy reported in initial experiments; secondly, their investigation reveals the limitations of relying on users' self-reported locations as the ground truth, this traditional approach consistently yields subpar outcomes (using sparser GPS-based locations as the ground truth, in contrast, yields considerably better results); thirdly, they shed light on the diminishing count of posts that a trained model can effectively tag, this reduction occurs at a rate of halving every four months, driven by changes in the platform's user base.

This thesis is particularly relevant because it tries to cover the whole subject of location inference of users in social media, while going extensively into details about location inference through social network-based models. Since many works cover the whole subject without trying to cover extensively the social network-based models, the peculiarity of this work is to cover in detail location inference through social network while still covering the various aspects of location inference, location data, and location privacy.





## 2 | Location Privacy Protection

Given the large amount of data present in social media today, in order to protect user's privacy, it is important that social media platform providers implement robust enough protection mechanisms. The abundant collection of data actuated by companies highlights the problem that these data need to be protected. The data collection still needs to continue for companies in order to run their businesses (usually advertising), but the privacy of the users need to be preserved as well. For this reason, there are privacy models and protection mechanisms to achieve privacy protection while preserving the utility of the data. To present to the reader the various methods that exists to achieve location privacy, this section will start by presenting the privacy models that exist. Afterwards it will describe the location privacy protection mechanisms that are used to achieve such privacy standards and properties. After that, it will cover the metrics used to measure the level of privacy protection. To conclude this chapter, last section will discuss the efficacy of the protection mechanism presented.

### 2.1. Location Privacy Protection Mechanisms

Numerous studies have dedicated their efforts to gathering relevant material on location privacy protection mechanisms [50][8][9][82]. Initially, privacy policy was guaranteed by a set of laws that ruled the relations between users and service providers. Such laws determined limits and conditions that external entities must follow when accessing, storing, and utilizing a user's location information in posses of the service provider. However, similar privacy policies are strictly dependent on the graphical locations of the service provider, which has to comply to the laws ruling that country. Essentially, the LBS (location-based service) server is obligated to obtain users' consent before accessing their location information. However, the LBS server not always completely respects such rules. For instance, the paper [32] examines 30 Android apps: about half of them, (the paper explicitly refers to MySpace and Evernote as the two of the most widespread platforms) leaked customers' location data to advertisement or analytics without explicit consent from the users. Also in case the LBS server strictly follows privacy rules, it is possi-

ble that it can be attacked by hackers and similar actors searching for information. For instance the article [56] quotes the case of the black hat hacker "Peace", who in 2013 exposed the data, including location information, of over 167 million LinkedIn users. A similar fate affected in 2016, about 360 millions MySpace users [77]. [48] adds the case of OpenSSL cryptography library which was attacked through a bug called Heartbleed. The hacker had a great possibility to extract sensitive data. The main problem connected to such system of privacy policy is that it relies on legal rules, but cannot be an aim in itself of the activity of the LBS provider. Moreover, such laws imply a follow of the frequent changes in rules and cannot be applied to dynamic location aware environment. To sum up, law-based privacy policies do not represent a sure and reliable protection of users data, they are always a step behind the technological developments and the growth of the amount of data made known in social media.

To protect users' privacy there are privacy models that defines properties that, when satisfied, will guarantee a certain degree of privacy. To first one is  $k$ -anonymity: The concept of  $k$ -anonymity, introduced in 2002 [103], aims at preventing the unique identification of individuals based on a small subset of their attributes known as a *quasi-identifier*. The *sensitive attributes*, which are not part of the *quasi-identifier*, are the attributes that need protection. The paper [82] considers the typical case of medical records. In this case, birth date, gender, and zip code are the *quasi-identifier*. Through *quasi-identifier* is possible to identify individuals. In this case the sensitive information is the disease of the patient. The principle of  $k$ -anonymity states that for effective protection, a user should be indistinguishable from at least  $k - 1$  other users. To this goal, all  $k$  indistinguishable users must share the same attribute values in their *quasi-identifier*, creating what is known as an *anonymity group*. Consequently, the likelihood of an attacker, without any external knowledge, to re-identify an individual among  $k$  similar users is limited to at most  $1/k$ . In the paper [82]  $k$ -anonymity is defined as follow: "Let  $d$  be a sequence of records with  $n$  attributes  $a_1, \dots, a_n$  and  $Qd = a_i, \dots, a_j \in a_1, \dots, a_n$  be the quasi-identifier associated with  $d$ . Let  $dk$  be the  $k$ -th record of  $d$  and  $r[Qd]$  the projection of record  $r \in d$  on  $Qd$ , i.e., the  $|Qd$ -tuple formed of values for only the attributes of  $Qd$  in  $r$ .  $d$  is said to satisfy  $k$ -anonymity if and only if each unique sequence of values in the quasi-identifier appears with at least  $k$  occurrences in  $d$ ". Formally:

$$\forall s \in \{r[Qd] | r \in d\}, |\{i \in N | d_i[Qd] = s\}| \geq k$$

In Table 2.1 there is an example of a dataset with  $k$ -anonymity with  $k = 2$ . The *quasi-identifier* are birth year, gender and zip code, and the *sensitive attribute* is the Disease.

This table provides 2-anonymity because by knowing birth year, gender and zip code from someone you cannot infer their disease, since there is at least another person with the exact same *quasi-identifier*. In the case where  $k = 3$  there would have been at least 2 other persons with the same *quasi-identifier*.

Birth	Gender	Zip	Disease
1980	M	0257	Asthma
1980	M	0257	Appendicitis
1980	F	0257	Neck pain
1980	F	0257	Asthma
1980	F	0257	Neck pain
1959	M	0222	High blood pressure
1959	M	0222	High blood pressure

Table 2.1: An example dataset with  $k$ -anonymity with  $k = 2$

$K$ -anonymity presents limitations, especially because it does not provide sufficient protection for users' location privacy. In case of densely populated areas, the region containing at least  $k$  users may be too small, which can lead to the disclosure of users' private data [20]. To extend  $k$ -anonymity,  $l$ -diversity is introduced [66]. It expands upon the concept of  $k$ -anonymity by introducing an additional requirement of having at least  $l$  well-represented values within each anonymity group. Specifically, it enforces a specific distribution of values for sensitive attributes among the members of each anonymity group. This concept of well-representation is formally defined in three different ways in [66]. The simplest form is referred to as distinct  $l$ -diversity, which stipulates that each sensitive field within an anonymity group must have at least  $l$  distinct values. The idea is to ensure that each sensitive attribute in the dataset has at least  $l$  distinct values within each group of records that share the same non-sensitive attributes. In simpler terms, it means that within a group of similar individuals (based on non-sensitive attributes), there should be enough diversity in their sensitive attributes to prevent the identification of a specific individual. For example, consider a dataset of medical records where each record includes attributes like birth year, gender, and medical condition. To achieve  $l$ -diversity, the sensitive attribute (medical condition) should be represented in such a way that each group of records sharing the same year of birth and gender contains at least  $l$  different medical conditions. This ensures that an adversary cannot easily identify a person basing on a unique combinations of attributes.  $t$ -closeness [57] is an additional enhancement to  $l$ -diversity, aiming to go beyond the sole requirement of having a diverse representation of sensitive values. In this approach,  $t$ -closeness ensures that the distribution of each sensitive attribute within anonymity groups aligns with the distribution of the attribute

in the entire dataset, with a threshold of  $t$  being applied.

Differential privacy is a model introduced in [28]: the concept revolves around ensuring that the computation of an aggregate result remains nearly unchanged, regardless of the presence or absence of any individual element within the dataset. Put simply, the probability of any outcome from an aggregate function should not be significantly altered by the addition or removal of a single element. Unlike  $k$ -anonymity, the definition of differential privacy remains unaffected by any external knowledge possessed by an attacker. Differential privacy is defined as: "Let  $\epsilon \in R^{+*}$  and  $K$  and  $K$  be a randomized function that takes a dataset as input. Let  $image(K)$  be the image of  $K$ .  $K$  gives  $\epsilon$ -differential privacy if for all datasets  $D_1$  and  $D_2$  differing on at most one element, and for all  $S \subseteq image(K)$ " [82],

$$\Pr[K(D1) \in S] \leq e^\epsilon \cdot \Pr[K(D2) \in S]$$

The authors from [5] introduce an expanded privacy concept called geo-indistinguishability, which extends the principles of differential privacy. They conclude that geo-indistinguishability can effectively safeguard the precise locations of individuals. According to this notion, the level of privacy protection desired should be linked to the concept of distance. Geo-indistinguishability can be defined as: "Assume a possible location set of users  $X$  and a probable reported location set  $Z$ , and let  $d(\cdot, \cdot)$  denote the Euclidean distance. For any two locations  $x_1, x_2 \in X, z \in Z$  and  $d(x_1, x_2) \leq r$ , algorithm  $K$  is  $\epsilon$ -Geo-differentially private if:

$$\Pr[K(x_1) = z] \leq e^{\epsilon \cdot d(x_1, x_2)} \cdot \Pr[K(x_2) = z]$$

where  $\epsilon$  indicates the privacy degree at one unit of distance" [50]. This definition implies that when the actual locations  $x_1$  and  $x_1$  are very close to each other, there is a similar probability in the distributions of generating the same new location  $z$ . On the other hand, as the distance between  $x_1$  and  $x_1$  increases, the difference between the probability distributions becomes larger. This difference is determined by the parameter  $\epsilon \cdot d(x_1, x_2)$ , which represents the level of privacy protection. It has to be noted that geo-indistinguishability consumes the privacy budget quickly, with the result of a finite number of LBS queries.

Now we can talk about Location Privacy Protection Mechanisms (LPPMs) to achieve privacy. The first one is *Generalization-Based Mechanisms*: Generalization techniques have been effectively utilized to achieve  $k$ -anonymity in the realm of location privacy, primarily applying "spatial cloaking" as proposed in [37]. The fundamental concept revolves around divulging less precise location data instead of the exact coordinates of the users. By doing so, it becomes possible to create cloaking regions where, at any given time, there are present at least  $k$  users. This not only diminishes the level of precision in the disclosed

information but also facilitates the establishment of areas where user identities remain indistinguishable. Usually there are three types of generalization: spatial generalization, temporal generalization, data generalization (also with fixed or adaptive ranges) [83]. A graphical representation of this mechanism is presented in Figure 2.1.

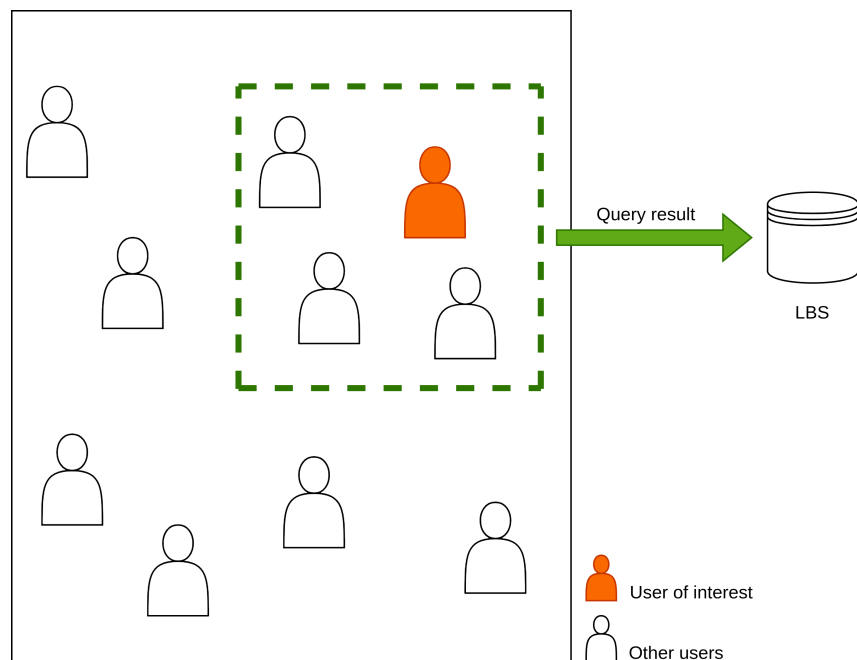


Figure 2.1: Generalization-Based Mechanisms.

*Dummies-Based Mechanisms:* another strategy to create  $k$ -anonymity consists in creating "dummy users" who are fake users. Such strategy works similarly to generalization based methods, but does not rely on the presence of other real users, which can hide their identity. Even if the attacker might be aware of the presence of fake users within the acquired data, he is not able to distinguish between real users and dummies. This method was first introduced by the work [54]. In general, dummy-based methods offer various advantages compared to other LPPMs:

- they operate independently of third parties;
- they enable accurate query results;
- they eliminate the need for sharing encryption keys between users and LBS servers;
- they maintain effectiveness even if attackers are aware of the privacy protection approach being employed.

However, these methods do come with certain drawbacks. Firstly, they tend to incur high communication costs. Additionally, they can result in resource wastage at LBS servers due to the servers having to respond to multiple fictitious queries. Lastly, only experiments can measure the real effect of these methods while mathematical proof in a strict sense can not be shown [50]. A graphical representation of this mechanism is presented in Figure 2.2.

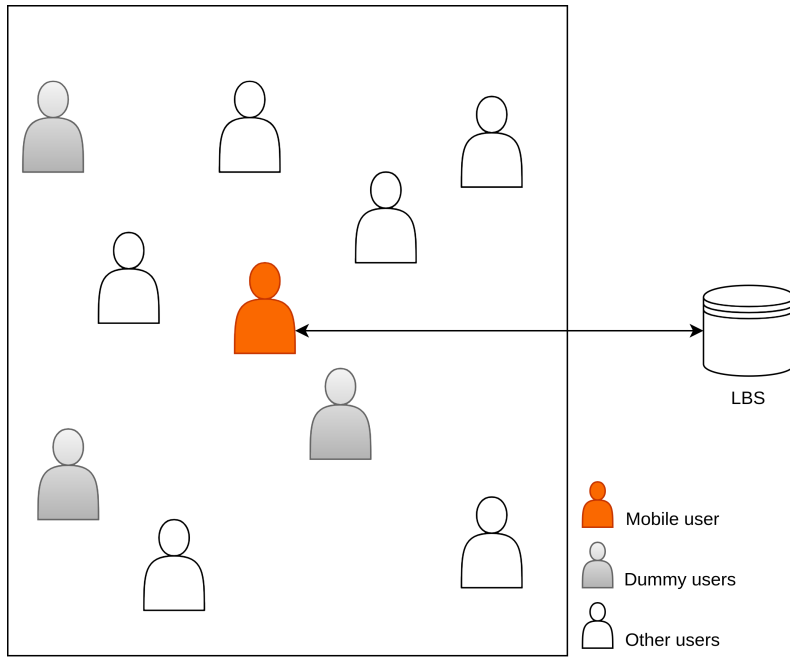


Figure 2.2: Dummies-Based Mechanisms.

*Hiding-Based Mechanisms:* there are two ways to apply these methods. The first one consists in the suppression of location data, the second one is their sampling. Noise introduction is avoided in these methods. According to [83] a possibility is to suppress small counts (SSC) by assigning zeros to location time-series which are lower than a predetermined threshold in the respective aggregate. Suppressing Less Popular Locations/Timeslots (SLP) consists in suppressing a percentage of the least popular visited locations and time-intervals. Sampling (SMP) consists in deleting a fixed percentage of users' data in a random manner, and in releasing the aggregates computed by utilizing the sampled data.

*Perturbation-Based Mechanisms:* protection techniques aim to achieve  $k$ -anonymity by concealing a user within a larger group; the alternative mechanisms aim at the same goal by modifying the data transmitted to an LBS for protection. Perturbation is usually used to achieve geo-indistinguishability. In this scenario, the challenge lies in finding a balance between privacy, where the data must be distorted sufficiently to safeguard it, and

utility, where excessively distorted data may render the results from the LBS unusable. Many mechanisms tackle this challenge by introducing additional noise, often in a random manner, to the original raw data. There are different types of perturbation: Perturbing Small Counts (PSC) involves introducing noise from the Laplace distribution to the counts of the aggregate location time-series that fall below a certain threshold  $k$ ; on the other hand, the Fourier Perturbation Algorithm (FPA) follows a specific procedure: firstly the Discrete Fourier Transformation changes the time series into the frequency; then the result is perturbed by adding noise using Laplace distribution. At last, perturbed time-series is obtained by the application of the inverse Discrete Fourier Transformation. Hiding and Perturbation can also be combined (Sampling Perturbing Small Counts, Sampling Fourier Perturbation Algorithm).

Another LPPMs is *Mix-Zones*: Mix-zones were first introduced in [144] after the pioneering work [145] on mix networks. The authors apply the mix-zones idea to the situation in which mobile users interact with LBSs: it is realized by using pseudonyms in the place of their real data that can unveil their actual identity (real name, IP, or MAC address). According to the authors, in a mix-zone the movements of users remain untrackable, because they cannot communicate directly with an LBS. Every time that users enters in another mix-zone, they are assigned with a new pseudonym. So, if there are  $k$  users at the same time in the same mix-zone, their identities will be mixed, achieving  $k$ -anonymity and creating confusion for potential attackers.

*Path Confusion* is another technique that mitigates the problem of the possibility of recovering trajectories from several location points, thanks to spatio-temporal correlations between them. To solve this with *Path Confusion*, when two entities meet, they exchange their identities with a specified probability so that the server cannot link consecutive location samples to one user or another.

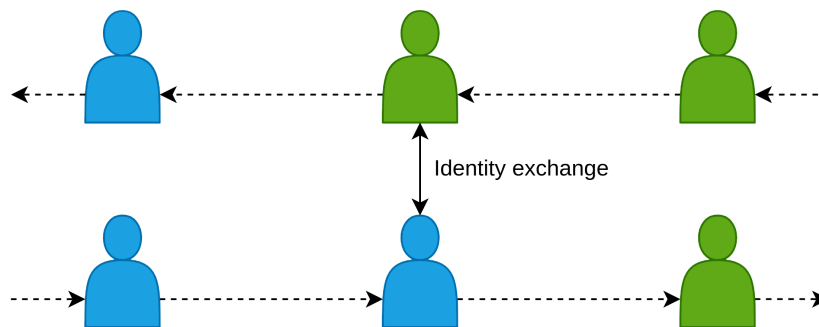


Figure 2.3: Path Confusion-Based Mechanisms.

Another LPPM is *Co-location Masking*: it consists in grouping together multiple nearby co-locations, limiting the ability for the attackers to identifying which of the co-locations in the group are real and which ones are false.

Lastly, there are techniques more focused on protecting users' privacy in graph data. Data like friendship, following/follower and mobility traces can be represented through graphs; for this reason techniques more oriented towards graphs anonymization are required. The first one is *Edge Manipulation*: an edge manipulation algorithm usually consists in randomly adding, deleting or switching edges to a graph structure.

Another technique is *Clustering* where an algorithm groups users and edges, and reveals only the information regarding the density and size of the cluster, protecting individual users data.

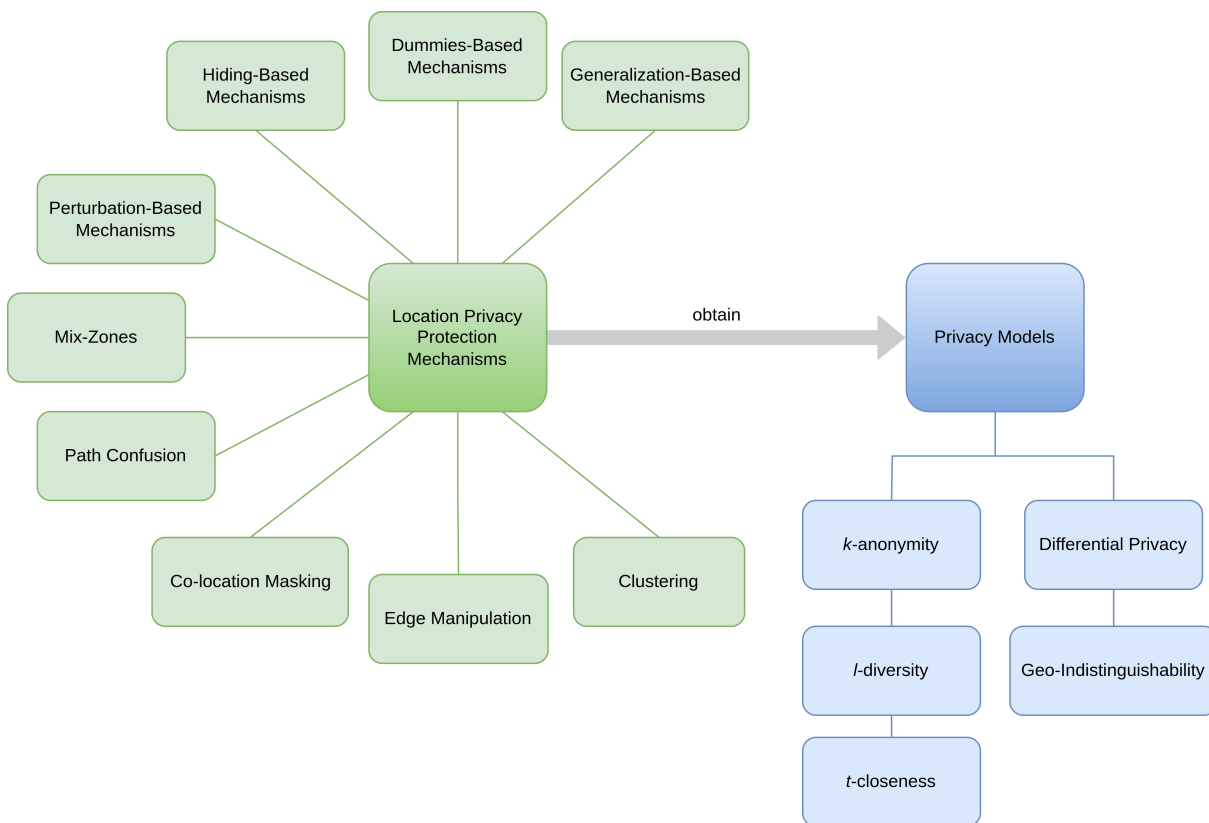


Figure 2.4: Conceptual map of LPPMs.



Name	Type	Explanation
$k$ -anonymity	Privacy model	A user should be indistinguishable from at least $k - 1$ other users. To this goal, all $k$ indistinguishable users must share the same attribute values in their <i>quasi-identifier</i> , creating what is known as an <i>anonymity group</i>
$l$ -diversity	Privacy model	Each sensitive field within an anonymity group must have at least $l$ distinct values
$t$ -closeness	Privacy model	The distribution of each sensitive attribute within anonymity groups aligns with the distribution of the attribute in the entire dataset, with a threshold of $t$ being applied
Differential Privacy	Privacy model	The computation of an aggregate result remains nearly unchanged, regardless of the presence or absence of any individual element within the dataset
Geo-indistinguishability	Privacy model	When the actual locations $x_1$ and $x_1$ are very close to each other, there is a similar probability in the distribution of generating the same new location $z$ . On the other hand, as the distance between $x_1$ and $x_1$ increases, the difference between the probability distributions becomes larger
Generalization-Based Mechanisms	LPPM	Divulging less precise location information instead of the exact coordinates of users
Dummies-Based Mechanisms	LPPM	Creating "dummy users" who are fake user. They are indistinguishable from real users
Hiding-Based Mechanisms	LPPM	Suppressing or sampling location data

Name	Type	Explanation
Perturbation-Based Mechanisms	LPPM	Introducing additional noise, often in a random manner, to the original raw data
Mix-Zones	LPPM	In a mix-zone the movements of users remain untrackable, because they cannot communicate directly with an LBS. Every time that users enters in another mix-zone, they are assigned with a new pseudonym. So the identities of users in a mix-zone are mixed
Path Confusion	LPPM	When two entities meet, they exchange their identities with a specified probability so that the server cannot link consecutive location samples to one user or another
Co-location Masking	LPPM	Group together multiple nearby co-locations, limiting the ability for the attackers to identifying which of the co-locations in the group are real and which ones are false
Edge Manipulation	LPPM	Randomly adding, deleting or switching edges to a graph structure
Clustering	LPPM	An algorithm groups users and edges, and reveals only the information regarding the density and size of the cluster, protecting individual users data

Table 2.2: Privacy models and LPPMs

## 2.2. Privacy evaluation metrics

Applying strong mitigation measures to aggregate data is technically possible, but it comes at the cost of reduced utility and increased overhead. When discussing defense mechanisms, a careful balance must be struck among three factors: location privacy,

utility, and overhead. Enhancing location privacy typically results in diminished utility and heightened overhead. With overhead we refer to the computational requirement increase to perform mitigation techniques. There are three types: computational overhead, communication overhead, storage overhead. Various metrics can be employed to gauge location privacy.

One is the anonymity level: The privacy concept seen before called  $k$ -anonymity requires the property of indistinguishability among a set of  $k$  users or dummy locations. The value  $k$  gives a representation of the privacy level, this is why we can talk about  $k$ -anonymity as a measure for location privacy, but also as a property for protection mechanism. So, for  $k$ -anonymity, the value of  $k$  is one of the parameters that can be used to represent the guaranteed location privacy. There is a problem with using  $k$ -anonymity: in the definition of  $k$ -anonymity there are no assumptions about the attackers' background knowledge. For this reason, the value of  $k$  fails to quantify the level of privacy in case of an inference attack with background knowledge.

Expected Estimation Error serves as the established criterion for quantifying location privacy. This measurement gauges the precision with which adversaries deduce the true position  $x$  by observing the disclosed location  $x'$  and utilizing their pre-existing knowledge. Normally the attacker's prior knowledge consists in the probability distribution across potential users' location. This privacy-preserving mechanism receives the genuine user location  $x$  as input and reveals a fabricated location  $x'$ . With  $x'$  in hand, potential attackers can calculate a posterior probability distribution pertaining to the genuine location  $x$  [50], denoted by:

$$Pr(x|x') = \frac{\pi(x)Pr(x'|x)}{\sum_{x \in X} Pr(x|x')}$$

According to [50]: "given the posterior probability distribution, adversaries can calculate an estimated location  $\hat{x}$  by minimizing the expected inference error, which is the expected deviation between the estimated position  $\hat{x}$  and the actual position  $x$ ":

$$\hat{x} = \underset{\hat{x}}{\operatorname{argmin}} \sum_{x \in X} Pr(x|x') |\hat{x} - x|$$

So, the expected estimation error is calculated as:

$$ExpErr_{privacy} = E|\hat{x} - x| = \sum_{\hat{x}} Pr(\hat{x}|x') |\hat{x} - x|$$

The effectiveness of various location privacy preserving mechanism, such as dummy locations, can be calculated through expected distance error.

Another metric is Entropy: it measures the level of privacy protection by assessing the uncertainty an attacker faces in identifying a specific answer from a pool of candidates. While this measure captures privacy in terms of the uncertainty an attacker has, it does not account for the accuracy of the attacker's estimation in relation to the actual location. In a similar way to the case of the expected estimation error, the attacker computes the posterior probability  $Pr(x|x')$  based on  $\hat{x}$ . Starting from the observed version  $x'$ , the attacker's uncertainty about the real location  $x$  can be represented as the entropy of the posterior probability according to the following formula:

$$H_x = -\sum_{\hat{x}} Pr(\hat{x}|x') \cdot \log(Pr(\hat{x}|x'))$$

The last measure is Geo-indistinguishability. As said before, geo-indistinguishability is a privacy concept that determines if an algorithm  $K$  is  $\epsilon$ -geo-differentially private. This measure is immune to side-channel information, the attacker's prior-knowledge, and implies that a user enjoys  $\epsilon r$ -privacy within  $r$  if any two locations at a distance of no more than  $r$  result in pseudo-locations with "similar" probability distributions. Therefore, parameter  $\epsilon$  represents the user's level of privacy. This parameter  $\epsilon$ , which quantifies the user's level of privacy, is the reason that geo-indistinguishability is also used as a measure to evaluate privacy. In fact this parameter, when geo-indistinguishability is used, essentially quantifies our level of privacy. As said before, geo-indistinguishability is a concept derived from differential privacy, wherein a method is often employed to introduce noise into the output of a query, which represents the pseudo-location. This noise adjustment alters the probability distribution of the pseudo-location, thereby achieving privacy protection. It is important to note that the privacy metric of geo-indistinguishability is primarily used to assess the privacy level provided by a location perturbation mechanism based on differential privacy.

On the other hand, when evaluating the utility aspect, a commonly employed utility metric is the expected distance between the released location and the actual location. This metric, often referred to as the "loss of quality", serves as a general measure of data utility. It is widely accepted by the research community and is applicable across various location-based applications, providing a standardized approach for assessing utility independent of specific application contexts. According to [50], the formula is:

$$ExpErr_{utility} = E|x' - x|$$

### 2.3. Effectiveness of LPPMs

The effectiveness of these systems when it comes to inferring location by using additional information from social media such as friendship network and co-locations can be questioned, on the basis of the various works that are presented later in this document. In the paper *Protecting Against Inference Attacks on Co-location Data* [2], it has been studied how effective are state of the art's defenses against attack that exploited users' co-location information. This was accomplished through the simulation of an adversary's attack, wherein the adversary aims to extract information from both the position of a user's check-ins and the positions of all other check-ins that possibly have co-located with the initial one. The objective is to identify patterns of user's check-in behavior and co-locations, following the idea that people usually co-locate with people they have social relationship with. The LPPMs studied are:

- Gaussian Perturbation: the location and timestamp of a co-location is distorted with a spatial noise vector and a temporal scalar noise magnitude derived from the Gaussian distribution.
- Adaptive Perturbation: it introduces noise dependent on the distribution of the nearest spatiotemporal neighbors for that co-location.
- Co-location Masking: it group together multiple nearby co-locations, limiting the ability for the attackers to identifying which of the co-locations in the group are real and which ones are false.
- $\epsilon$ -geo-indistinguishability: they considered the mechanism implemented in [5] for achieving geo-indistinguishability.

From their experiments it results that:

- Gaussian Perturbation: as noise increases the attack precision decreases, but the protection is inconsistent. In fact, users in sparse areas are still not well protected.
- Adaptive Perturbation: the Adaptive Perturbation mechanism solves the problem present with Gaussian Perturbation, where sparse areas are no well protected. With adaptive perturbation that does not happen anymore. As protection increases, the inference precision decreases.
- Co-location masking: the obtained results are similar to the one obtained with Adaptive Perturbation.
- $\epsilon$ -geo-indistinguishability: the precision of the attack results lower, however it re-

sults that the protection is inconsistent. Again, the areas that are less protected are the less dense ones. It also results that with some implementations of geo-indistinguishability the noise introduced is so large that it makes the data useless.

Their experiments make evident that state-of-the-art methods not tuned specifically against co-locations introduce too much noise in the data. This results in the data being not very usable; more importantly, the protection against attacks that exploit co-location information yields poor results. In conclusion, further research is required to advance the existing state of location privacy protection mechanisms in addressing location inference through co-location and social relationships.

## 3 | Location Inference

Location inference is a type of attribute inference, as defined above. As said in [79]: "attribute inference can be regarded as a method to infer a group of sensitive attributes that users don't want to be known by others, from users' online publishing and interaction information". So, location inference is the practice of inferring the location of users without them disclosing it directly. This can be achieved by using location related information contained in users' post, such as mentioned locations or word used mainly in specific geographic regions. Also user friendships and social relationships like co-workers and such can be leveraged to infer someone's location more accurately. In this section various works that focus on inferring location of user's social media are collected. Sections 3.1, 3.2 and 3.3 will refer to the various types of locations that can be inferred respectively: home location, which is a user's residence location; activity location, where users tends to spend time outside of their house; and next location, which is a future place that a user might go to in the future. Then, Section 3.4 will cover the most important part of this thesis, which is location inference through social network. By social network, this work refers to the various social relationships that a user can have with other users. This section is divided in two subsection: Subsection 3.4.1 will cover works that utilize information from friends of a user, Subsection 3.4.2 will cover works that utilize information from other users that have similar living patterns to the user of interest. This division has the following reason: we can learn important information from friendships that show themselves on online social media in an explicit manner via mentions, interactions, followings and so on. But some friendships do not manifest online and there is no way to find those real life friends other than by observing the similarity in their moving patterns. This because friends tends to visit places together, do the same activities and share some similar living patterns. In the process of analyzing similar users, it is also possible to retrieve a lot of information from users that have high similarities, but are not necessarily friends, since they share similar living patterns.

### 3.1. Home Location Inference

Inferring the home locations of social media users has numerous applications, including local content recommendation, location-based advertising, public health monitoring, and public opinion polling, etc. However, since it is often not mandatory for users to provide their residential information on most social media platforms, their home locations are frequently missing or unreliable. While location can still be derived from a user’s profile, this approach has limited effectiveness, particularly on platforms like Twitter where the information is typically limited to city or state level [102]. As a result, considerable research has been dedicated to deducing the home locations of users.

When discussing posts composed from text (like tweets for example), the primary source of information we typically rely on is the text itself. An effective way to infer location from tweets’ text is to identify location relevant words [121]. Most studies model the problem of using local words by using a probabilistic model. A representative probabilistic model is introduced by the authors of [15], as follows: "the distribution of users’  $u$ ’s home location  $l$  given their tweet posts  $S(u)$  is decomposed as  $P(l|u) \propto \sum_{w \in S(u)} P(l|w)P(w)$ " [121]. They considered local words  $w$ , and  $P(w)$  stands for the probability of  $w$  over the entire corpus. The focus is to estimate the *spatial word usage*, which consists in the location distribution  $P(l|w)$  of word  $w$ . With this strategy, they were able to estimate a user’s city-level location, even when the posts did not have any location cues.

An effective way to infer locations is to analyze the context of a post, including location description and timezone. In [88] the authors used a multilayer perceptron (MLP) with one hidden layer in order to infer users’ home locations. They employ the l2 normalized bag-of-words representation of a user’s tweet content as input, with the outcome being a pre-defined discretized region formed using either a k-d tree or k-means algorithm. The authors from [67] and [68] used as source of information to infer locations tweet posting time. Posting times are logged in GMT (Greenwich Mean Time). They divide the GMT day into slots of equal-length times, so that they can view users according their distributions. Due to the variations in time zones, users display shifts in their distribution patterns. To address this, a time-zone classifier is trained using the distribution as input feature. These classifications offer insights into the time zones of users and can provide a broad estimation of their locations. Since self-declared locations and time zones in posts are often deceptive because for instance abbreviations can be used to indicate places, the authors of [40] and [41], they also incorporate features such as four-grams of self-declared locations and time zones to train a classifier for identifying home locations. In [29] the authors developed a probabilistic model that recourse to the temporal distribution of



geo-tags contained in tweets in order to estimate a user's home location and workplace. They ground their method on the assumption that there is an high probability that late in the night users' tweets come from their home location, while in day time hours they come from what could be their workplace. The authors from [81] affirm that relying only on the media as home locations could be deceptive, because users often can have multiple active regions. To avoid this error, they consider as home clusters the group of posts with the highest recurrence after having clustered all the geo-tags. By making the geometric median of all points within the home clusters, they assume to attain the home coordinates. Another possibility, considered in [16], is offered by the aggregation of users' geo-tags into a square pattern. By pinpointing the most geotagged square they inferred the central point, and by iterating this procedure within the central area, they break it into smaller unities until they reach pre-established threshold, so they do not simply apply the geometric median, so that the user's home location corresponds to the final center of the procedure. The work [87] takes a different approach by leveraging a neural network model combined with a mixture density network. They transform the two-dimensional geo-tags into a continuous vector space and use them as input for inference. An alternative approach to deducing a user's home location entails examining the distribution of their check-ins and photos, as exemplified in [100].

## 3.2. Activity Location Inference

Instead of inferring users' residence location, there may be an interest in determining the location they visit when they are not necessarily at home. This can be accomplished by inferring the location from where a user makes a post on social media platforms like Twitter.

To infer the location from tweets, word-centric methods can be employed [121]. In [1] the authors propose an approach that utilizes Gaussian mixture models for tweet location prediction. They focus not only on modeling the spatial usage of individual words but also n-grams to increase redundancy since only one tweet is available as input for tweet location prediction. Another similar approach is to model spatial n-gram usage using Gaussian models, as discussed in [33]. The authors of [18] use a different approach, based on a learning-to-rank method: to encode tweet content it recourse to smoothed probability estimation of words recurring at a specific venue.

Instead of focusing solely on words and n-grams, another approach is to prioritize the topic of the tweet. This can be achieved by extracting the topic from the words used in the tweet. Words such as "NBA" and "Kobe" could represent the topic of "basketball", as

exemplified in [121]. To better illustrate the work made on topics it's better to talk about two works that focus on home location inference, but are used as reference by successive works for location prediction by focusing on topic. The work [30] introduces the concept of corrupting topics by including related terms. In the example taken the corrupted topic for "basketball" may also include "Celtics" corresponding to the location in Boston, which is the city of that team. The same authors further enhance the concept of topic corruption in [31] by formulating what they call Sparse Additive Generative model (SAGE). This model develops the notion of location-based topic corruption already presented in [30], while also enabling sparsity and simplicity in model inference. However, the reader of this paper must keep in mind that the works quoted [30] [31] are reported here on for sake of completeness about topic-model-based methods, because they focus on home location inference only. The SAGE model [31], has been extended by other works, as for instance [44]. This paper considers a set of data (region, topic and users' interests) to build a model that maintains the original structure of tweets, differently from [30] and [31]. In this way to model location in a per-tweet manner, the authors assume that tweet location depends on the user's geographical interest distribution, while the topic of a tweet relies on the user's topical interest as well as local topics. In this perspective, tweets are created starting from the topic and according to a local word distribution. Paper [14] assume another point of view. It does not model users geographical interests in function of a multinomial distribution. Instead it uses users' interests as a latent variable to construct a specific function related to the location (eating, shopping, health care and similar). All this variables, are used to bridge users and locations since each user is characterized by a specific distribution interest according to location and related activities, which are involved in tweets generation. Work [117] presents another approach, which show some similarities to the previous one, but presuppose the recourse to an intermediate variable: the authors call it regions and it is located between users and tweet locations. This by following the idea that users normally have a "work" place (region) and a "home" place (region), which are modeled by Gaussian distributions having center coinciding with the respective workplace and home location. Following the supposition that users are involved in activities near their work and home places, for example during lunch time, they would often choose a restaurant near their work place. By doing so, it is probable that the user would leave a tweet about eating there, quoting the restaurant's name.

Contextual information, such as location description and timezone, can also be leveraged [121]. In fact, a timestamp may be very informative if enough historical data for locations are provided; for example, bar, public local and clubs are normally posted about during afternoon and night, according to normal working time. Instead, places such as parks

are tweeted about during weekend days or festive days. By keeping this in mind, the authors from [60] divide tweets according to time distributions for locations at three different scales: day, week, and month. Preferences between location are inferred by the combination of tweets associated with the same timestamp, according to the set of three distributions. In the same way but on a lower scale, the work [117] recourses to a binary distinction between weekdays and weekends combined with the different time in a day. So, by basing on the users choices, the authors can infer if a tweet is posted during a weekday or weekend. Given a user, the generative model determines whether a tweet is sent on a weekday or weekend. In a successive step they can infer also the daytime again from their daily hours distribution. Again, time zone and tweet posting time are both taken into consideration in [27] to use them as features in a classifier. The scope is to predict users' location according to cyclical temporal patterns as reconstructed by the classifier. In a more traditional way the authors of [99] collect information from data made accessible by the users in their public profiles such as self-declared home locations, websites, and timezones, and so on. All that in order to create what they call "polygon-shaped administrative regions" created by a query among many different databases for 10 specific indicators. They assume that the height of the polygons corresponds to these 10 indicators, so that in a following paper they can infer possible tweet locations [110]. Experiments confirm that such a multi indicator approach is more effective than single indicator approach, because the latter seems to be affected by ambiguity. From another point of view, paper [18] is based on using context information by assuming that location prediction can be inferred with great probability from both venues active time and users' visiting place histories. In first place, they use a smoothed kernel density estimation method to evaluate the possibility that a user can get to a location at a time given by analyzing venues active time. After that, they conclude that a normal user is spatially is usually limited in spatial choices because of a set of motivation (geography, social conditions and relationships, personal factors). So, by combining the previous information and considerations they presume to better estimate if a user could be in a specific location.

Other work demonstrate that are other ways to infer location that use less conventional sources to gather information to infer activity location. For example the work *Tagvisor: A Privacy Advisor for Sharing Hashtags* [118] presents the first systematic analysis of privacy issues related to hashtags, which have not been explored before. By using a random forest model, they showed that they can infer a user's precise location from hashtags with accuracy from 70% to 76%. In the work *PowerSpy: Location Tracking Using Mobile Device Power Analysis* [71] they take advantage from the fact that modern mobile operating systems such as Android allow installed applications to read aggregate

power usage on the phone. This information is considered harmless in regards of privacy, but with the use of machine learning algorithms, they managed to infer location using that information. To locate the phone they need strong assumptions about the attacker prior knowledge: they assume the attacker has prior knowledge of the area or routes through which the victim is traveling. This knowledge allows the attacker to measure the power consumption profile along different routes in that area in advance; they also assume that the tracked phone is moving by car or by bus while being tracked, because their system cannot locate a phone that is standing still since that provides only the power profile for a single location. Their experiments show that it is possible to track users who follow a daily routine. For example a mobile device owner might choose one of a small number of routes to get from home to work; the system identifies what route was chosen and in real-time can identify where the phone is along that route. This work is important because it underlines also how some apparently harmless information can be privacy inferring. Another unconventional way to infer location has been found in the paper *Inferring User Routes and Locations Using Zero-Permission Mobile Sensors* [73], where the authors discovered that they can use the gyroscope, accelerometer and magnetometer data on the smartphone of online social network users to infer the user's location. This can be done by modeling the problem as a maximum likelihood route identification on a graph. The graph is generated from the OpenStreetMap publicly available database of roads. From their experiments, it results that for most cities it is possible to output a list of 10 routes containing the traveled one with a probability higher than 50%.

### 3.3. Next Location Prediction

Instead on inferring the home location of a user, companies can be interested in predicting the next location of a user to improve things such as advertising etc. Usually, it is easier to do that with data coming from LBSNs. The mobile application of Foursquare and of other LBSNs allows us to know the exact location of a user at a given time through users' check-ins when they visit places. Location prediction involves training a model using users' historical tracks to anticipate their next position. Typically, significant places are extracted from the trajectory history, and a statistical model is employed to predict the user's next location [35].

The first category is Machine learning-based prediction. The study *MPE: A mobility pattern embedding model for predicting next locations* [12] introduces a new mobility pattern embedding a model called MPE, which exhibits two distinctive features. The first feature is integration of different types of information, mainly locations, time and object: they

are combined into a low-dimensional latent space, which is an abstract multi-dimensional space which encompasses feature values that are difficult to interpret directly. The other feature is the evaluation of the data originating from road networks in traffic trajectories. The proposed method surpasses state of the art methods, according to the experiments conducted by the authors. Another work *Exploiting machine learning techniques for location recognition and prediction with smartphone logs* [17] develops a method that is grounded on Hidden Markov models: they utilize k-nearest neighbor and decision trees, for inferring the users' next locations, gaining 90% of successful results. In the study *Unlicensed taxis detection service based on large-scale vehicles mobility data* [106] the authors aim to detect unlicensed taxis in a context of huge traffic by observing personal mobility trends. In a first place a set of spatio-temporal data are taken from the information collected, and second a system of three techniques of machine learning (SVMs, decision trees and logistic regression) are used to correct accuracy of information. This study is conducted in Xiamen, China, with excellent result, even though some adjustments are related to other datasets (demographics and POIs and so on). The work *Where are you going? Next place prediction from Twitter* [21] proposes two models: one for exploring individual features and a related one for combining them into a supervised learning framework grounded on a M5 decision trees incorporated in a WEKA framework (machine learning algorithms). The list of considered features includes: weekend location, visit count, visit probability, nighttime location, tweet count, location visit entropy, movement entropy, popular destinations, and spatial characteristics. They have taken into account a multitude of users' features. The algorithm chosen is the NexT clustering algorithm which gives an accuracy of 83%.

The second category is Deep learning-based prediction. In the study *Predicting the next location: A recurrent model with spatial and temporal contexts* [61], a Spatial-Temporal Recurrent Neural Network (ST-RNN) is proposed. This model incorporates local temporal and spatial contexts in each layer to uncover mobility patterns. The results achieved in this work are superior to existing state of the art's approaches. Another work *Location embedding and deep convolutional neural networks for next location prediction* [97] deals with the training of a model which classifies the sequence of previous location corresponding to a user in order to infer successive locations. On this basis, the authors incorporated a model called loc2vec on the basis of the similar model word2vec, which is used to predict the successions of words in a phrase. The new system encodes locations into a vector and when two or more location are present in the sequences at the same time, the vectors will be at close distance. This way, this method improved the state of the art. The research paper *CEM: a convolutional embedding model for predicting next locations* [13]

introduces the CEM model, which aims to predict future locations using traffic trajectory data. The model utilizes a one-dimensional convolution to capture the relative ordering of locations, taking into account the constraints imposed by road networks. Additionally, CEM incorporates a double-prototype representation for each location to eliminate incorrect location transitions and considers a combination of factors (such as sequential, personal, and temporal) that influence human mobility patterns. By encompassing these elements, CEM achieves higher prediction accuracy compared to methods solely based on sequential patterns. The effectiveness of CEM is validated through experiments conducted on two real-world trajectory datasets, where it outperforms existing state-of-the-art approaches. In the work *T-CONV: a convolutional neural network for multi-scale taxi trajectory prediction* [65] the focus is on predicting taxi trajectories to enhance various intelligent location-based services, particularly targeted advertising for passengers. The paper introduces T-CONV, a model that represents trajectories as two-dimensional images and leverages multi-layer convolutional neural networks to capture multi-scale trajectory patterns. Through comprehensive experiments using trajectory data, T-CONV achieve results superior to existing state of the art’s approaches.

Works	Inferred Location	Data	Method
[15]	Home location	Location relevant words	Probabilistic model
[88]	Home location	Posts’ context	Multilayer perceptron
[67]	Home location	Posts’ context	Time-zone classifier
[68]	Home location	Posts’ context	Time-zone classifier
[40]	Home location	Posts’ context	Self-declared location and time-zone classifier
[41]	Home location	Posts’ context	Self-declared location and time-zone classifier
[29]	Home location	Posts’ context	Probabilistic Model
[81]	Home location	Posts’ context	Geometric median of locations
[16]	Home location	Posts’ context	Geometric median of locations
[87]	Home location	Posts’ context	Neural network and mixture density network
[100]	Home location	Posts’ context	Probabilistic model

Works	Inferred Location	Data	Method
[30]	Home location	Posts' topic	Topic corruption
[31]	Home location	Posts' topic	Topic corruption and Sparse Additive Generative model
[1]	Activity location	Location relevant words	Gaussian mixture model
[33]	Activity location	Location relevant words	Gaussian model
[18]	Activity location	Location relevant words and posts' context	Learning-to-rank
[44]	Activity location	Posts' topic	Sparse Additive Generative model
[14]	Activity location	Posts' topic	Bayesian model
[117]	Activity location	Posts' topic and context	Probabilistic generative model
[60]	Activity location	Posts' context	Language modeling
[27]	Activity location	Posts' context	Feature classifier
[99]	Activity location	Posts' context	Multi indicator approach
[110]	Activity location	Posts' context	Multi indicator approach
[118]	Activity location	Posts' hashtags	Random forest
[71]	Activity location	Phone power usage	Hidden Markov model
[73]	Activity location	Phone gyroscope, accelerometer and magnetometer data	Maximum likelihood route identification
[12]	Next location	Time and locations	Low-dimensional latent space
[17]	Next location	Locations of nearest neighbors	Hidden Markov model
[106]	Next location	Taxis' mobility	SVMs, decision trees and logistic regression

Works	Inferred Location	Data	Method
[21]	Next location	Tweet count, visits count, nighttime location, weekend location, popular destinations and spatial characteristics	WEKA framework
[61]	Next location	Local temporal and spatial contexts	Neural Network model
[97]	Next location	Location sequences	Deep learning based on vectors
[13]	Next location	Traffic trajectories	Double-prototype representation
[65]	Next location	Taxis' trajectories	Multi-layer convolutional neural network

Table 3.1: Works regarding various types of Localization through Social Media

### 3.4. Location Inference through Social Network

Rather than solely concentrating on individual users to deduce their location, incorporating supplementary data from other users has demonstrated substantial enhancements in localization attacks. Real-world relationships frequently extend to social media platforms, as users commonly engage with those in close proximity to their residential areas. This reality underscores the significance of social relationships as an attribute that can significantly enhance the precision of inferring an individual's location.

Furthermore, an alternative approach to harnessing the collective insight of other users for location inference involves analyzing clusters of analogous users. This analysis aims to discern patterns and various forms of resemblance related to location, encompassing factors like residence and check-ins.

#### 3.4.1. Location Inference through Friendships

A common way to infer location of users using friends location data is using a probabilistic method that focuses on maximizing the likelihood of a location to be the user's



location. The work *Find me if you can: improving geographical prediction with social and spatial proximity* [7] proposes a probabilistic method based on the distance between users in the friendship graph. The FindMe algorithm makes the assumption that geographical distance of two users is inversely proportional to the possibility that they be friends in the network. So, starting from the distribution of friends locations, the method chooses the location that maximizes the likelihood of it being the user's location. This method has been reworked by *SPOT: locating social media users based on social network context* [55], where they assume that connections of a users should have different weights when inferring user's location. On this assumption, there is the need to utilize a coefficient for social proximity, which represents the similarity of the friendship network of two users. The number of common friends in the social network between two users evaluate the similarity of their networks. The coefficient so calculated can help to determine the weight of the location of a friend in the case of location inference. Another work *Location prediction in social media based on tie strength* [70] extends the method from [7] adapting it from Facebook to Twitter. The focus of this work is to classify a user's Twitter relationships because it probably serves as predictors of a user's location. After having selected users' with a minimum of three GPS posts, and calculated the average latitude and longitude values to establish their location, a tree regression model is trained with the established relationships between them, so that it became possible to predict their distance. In this way the predictive capability are determined by the predicted distances. The social relations are divided into ten parts, ordering them from the most predictive to the least predictive. In this way we can identify the friendships that are most useful for location inference. During the likelihood calculation, the most useful friendships found in the previous step will carry more weight. These relationship partitions are what separates this method from the one of [7]. The last work that this thesis presents that extends the work [7] is the work [69]. They investigated the relationship between the distance between a pair of two users, and the strength of the tie between them. They used several factors for their work such as: following, mentioning and conversations.

In another paper, *Towards social user profiling: unified and discriminative influence model for inferring home locations* [59], the authors aim to deduce geolocations by considering the combined impact of both users and locations. This approach stems from the notion that specific users hold more valuable insights for predicting the locations of their neighbors. The most effective strategy among their approaches initially assigns users to random locations, subsequently updating these locations in an iterative manner based on information from neighboring users and mentioned location references. This refinement process involves adjusting parameters in response to the prediction error of users with known locations.

Lastly, the work [84] uses a probabilistic framework based on a semi-supervised factor graph model to infer location of users. To better improve accuracy, they exploited the users' network made of friends from social media.

Bayesian network (a type of probabilistic graphical model) is a common method to infer location of users. In a study [74], titled *Quantifying interdependent privacy risks with location data*, the authors propose an approach that infers a user's location from reported co-locations in messages and pictures posted on social media. Another source of co-locations is retrieved from IP addresses of user's friends. These co-location information is also extended to friends of friends and beyond. A general Bayesian network model is the basis for a belief propagation algorithm proposed in the study for inferring locations. In the work *Finding your friends and following them to where you are* [96], they studied the correlations between trajectories of a user's friends. To do that they used a Dynamic Bayesian Network trained of locations sequences. The work [75] implements a Bayesian hidden location inference model and a multi-factor fusion based hidden location inference model that also exploits friends' locations. The addition of this paper is that it measures the similarity between two friends. This concept is extended in other works that we will see later when this similarity is not only used for friends. In the study *Multiple Location Profiling for Users and Relationships from Social Network and Content* [58] the authors explore the concept that users could potentially have multiple home locations. To uncover these multiple home locations, they employ a sophisticated approach that leverages social connections between users and references to location names found in posts. This is achieved through an extended version of the supervised Latent Dirichlet Allocation, a type of Bayesian Network.

Artificial Neural Networks are a common tool in this field. In the work '*Current city' prediction for coarse location based applications on Facebook* [10] the focus is directed to extracting explicit and implicit location information from users' friends. This information is then utilized to build a Current City Prediction Model (CCP) using an artificial neural network (ANN) learning framework. This model aims to accurately predict the current city of target users. Another work that uses a neural network is [72]. In this work they encoded user friendships into a neural network model. This neural network model learns from unified text, metadata and from the user network that is generated from user mentions.

Another work is [101], that introduces a de-anonymization attack that uses a user's friendship information in social media to de-anonymize users mobility traces. The underlying concept is that people meet with those who have a relationship with and so they can be identified by their social relationships. A contact graphs is created through the mobil-

ity traces of users by using Distance Vector, Randomized Spanning Trees, and Recursive Subgraph Matching heuristics. These methods permit to calculate the mapping strength and propagate it through the network. The accuracy and computational complexity of this work was later improved in [49].

There are many works that focus to achieve accurate location inference basing their methods on the label propagation algorithm (LPA or LP algorithm) [124]: the LP algorithm can be defined as: "a semi-supervised, iterative algorithm designed to infer labels for items connected in a network". The LP algorithm is by far one of the most used methods in this field. In Figure 3.1 an example illustrating how the algorithm works is presented. It has to be noted that, since the LP algorithm is not deterministic, the last node in the graphs can obtain the blue label or the orange one. In some variations a node can obtain two labels at the same time.

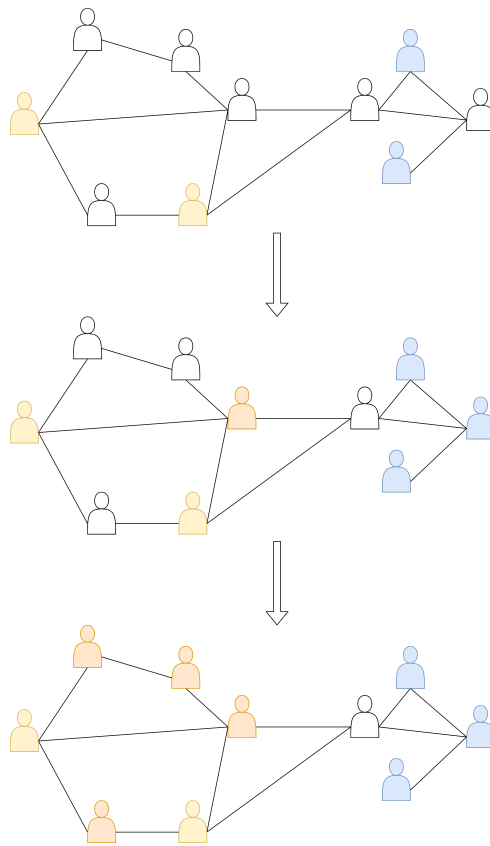


Figure 3.1: Label Propagation Algorithm example.

The first work that presented in this thesis that utilizes the LP alg. is called *That's What Friends Are For: Inferring Location in Online Social Media Platforms Based on Social Relationships* [51]. The authors allocate a user's location as the geometric median

of all the locations of their neighboring users, utilizing ground truth data. Subsequently, these initial assigned labels can serve as a basis for deducing a user’s location during a second iteration through the dataset. The research presented in *Geotagging one hundred million Twitter accounts with total variation minimization* [23] addresses the challenge of deducing user’s locations as an optimization over a social network with a total variation-based objective. The study introduces a scalable and distributed algorithm to tackle this problem, building upon the adaptation of the label propagation method outlined in a previous work [51]. This advancement in [23] refines method from [51] by enhancing the initial location calculation process in the first pass. This enhancement incorporates weighted edges in the network, where the weights are determined by the frequency of interactions between user pairs. This adjustment favors locations of friends with whom a user interacts most frequently. Subsequent iterations of the label propagation algorithm focus on updating a user’s location only if it aligns reasonably closely with the locations of their neighbors. In this way, errors in location information cannot be propagated. In the work *Find you from your friends: graph-based residence location prediction for users in social media* [113] the authors adapt the LP algorithm with the assumption that users geographically close to each other are more likely to establish friendship relations. They also pose that nearby users are inclined to share more common friends and publish similar geo-related content. This assumption stems from the notion that users’ social media content is primarily influenced by their physical world experiences. This approach surpasses the FindMe method in terms of accuracy and efficiency. While state-of-the-art techniques necessitate greater distance between non-friends, this method leverages the presence of more friends in close proximity. Its scalability and efficiency are enhanced by considering only a user’s friends instead of the entire user population, thereby reducing the dataset volume. With location inference we can focus on inferring the residence location of users, or the locations where a user tends to spend time. These types of location are referred in [39] as “activity locations” and are defined as the places among the highly checked-in region where most of user’s activities occur. Inferring these types of locations is the focus of the work *Activity location inference of users based on social relationship* [39], where they exploit social relationships following the assumption that users’ location tends to coincide with the locations their friends go to with high frequency. They proposed a method where they ignore friends whose majority of the neighbors have far activity locations. First they begin by verifying a user’s location based on the distance-probability relationship for friendship. To do this, they simultaneously select each labeled user to validate whether the location information of the users can be exploited to predict the locations of their neighbors. Therefore, after selecting a labeled user, they mask the user’s location first. After that, basing on the location distribution of the user’s

labeled neighbors, the location among the neighbors' activity location that has maximum likelihood is assigned. In case where the last assigned location of a determined user appears to far from the original activity location, this user will be removed from the pool of users taken into consideration for location inference. After addressing the neighbor's selection challenge, the study delves into the core issue of label propagation, that is how to prioritize inferring locations for unlabeled users. The priority is set as follows:

- users having a number of labeled neighbors closer to their mean location point;
- users with number of labeled neighbors;
- users with number of neighbors (both labeled and unlabeled).

These represent the fundamental concepts underpinning the sequential iterative algorithm introduced in this research, termed as Sequential Spatial Label Propagation (SSLP). The evaluation of this approach involved the utilization of several key parameters: average error distance, accuracy, inference coverage (a metric gauging the proportion of unlabeled users within a dataset that received a location assignment irrespective of accuracy), and running time. From testing, SSLP results to have lower average error distance than state of the art methods such as [7] and [51]. Same thing happens for accuracy where SSLP beats state of the art methods. For users inference coverage they noticed that SSLP has little lower coverage than SSLPn, which is the same algorithm as SSLP without the neighbor validation process. When it comes to running time, SSLP happens to perform worse than [7] and [51] methods by a significant amount. As expected, SSLPn performs worse than SSLP. Alternative approaches to employing label propagation for location inference have been explored, as exemplified in the study *Twitter User Location Inference Based on Representation Learning and Label Propagation* [104] where they use representation learning and label propagation (ReLP). Their method involves a series of steps: connection relation graph construction, user relationship filtering, user representation learning, propagation probability calculation, and user's location inference. To initiate the construction of the connection graph, user's relationships are established by utilizing regular expressions to identify mentions of "@user" within user-generated texts. Additionally, the information gain rate (IGR) of all words is computed, with words exhibiting high IGR considered indicative of location. These location-indicative words, combined with collected user mentions, contribute to the creation of an undirected connection graph. In the process of filtering user's relationships, an algorithm is implemented to eliminate predominantly global celebrities who are frequently mentioned but lack relevance to user-geographic attributes. In the third step (user's representation learning), they map location-similar users to a vector space based on a connection relation graph. For

propagation probability calculation, two key phases are executed. Initially, directly or indirectly adjacent users are identified through the adjacent relations of the connection relation graph. Subsequently, a propagation probability matrix is established using the user's relationship matrix. Notably, the likelihood of label propagation between users is heightened when their geographic attributes exhibit greater similarity. Finally, for the user's location inference step, the label propagation algorithm is actuated. Notably, the study acknowledges the challenge posed by users maintaining distinct accounts across various social network platforms. Consequently, the conventional practice of predicting a user's location solely based on a single social network is deemed limiting in this context. Another work that uses the LP algorithm is *Twitter user geolocation using a unified text and network prediction model* [85]. The data that they used for analyzing the social network are @-mentions. The same authors published also another work [86] that shows that Label propagation is very powerful when used in conjunction to text-based methods. Similarly to the previous one, this work uses an hybrid method focusing on @-mentions between users.

In the research paper *Where's @wally?: a classification approach to geolocating users based on their social ties* [94], they model the problem within the classification framework: a user could be assigned to any of the cities in the United Kingdom. To achieve this, a Support Vector Machine (SVM) classifier is utilized, incorporating three primary features: (i) the cities of the user's friends, relative to the number of Twitter users in that city, (ii) the number of closed triads in a user's social network residing in the same city, and (iii) the number of reciprocal following relationships a user has per city. A simple method for location inference is the one proposed in the work *Inferring the Location of Twitter Messages Based on User Relationships* [26]. In this work the Neighborhood Vote method is proposed. This method makes the assumption that the location of a user is the same as the most frequent location occurring within its friendship group. In order to correct the limitation implicit in this model, determined by whether a user has too many or, on the contrary, too few connections, they introduce a series of corrective parameters, which must obviously be adjusted: (i) minimum number of connections, (ii) maximum number of connections, (iii) minimum frequency of mentions guaranteeing the reliability of the inferred location. The paper *Toponym disambiguation in online social network profiles* [34] utilizes a social graph along with friends' locations. It introduces LocusRank, an algorithm designed for location inference within online social networks. LocusRank constructs a social graph utilizing the self-reported (potentially ambiguous) locations of a user's friends. This method then disambiguates the target geographic location and accurately infers the user's likely location with a high degree of probability.

The work [11] focuses on inferring home location for Twitter users. To solve this problem, they proposed a Social Tie Factor Graph Model (STFGM) for inferring users' city-level home location using three factors: (i) following network, (ii) user-centric data and (iii) tie strength.

In the work [38] they introduced a new concept called social trust. This social trust concept quantifies the number of mutual friends, and by doing so it measures the closeness of two entities. This is obviously leveraged to better infer locations.

The work [114] tries to solve the problem that some friends live apart from each other. To solve that, they introduce the concept of landmarks: users with a lot of friends who live in a small region, and not far from each other. On this basis they succeeded in inferring users' location through a model introduced by them, called Landmark Mixture model.

A possibility to improve location inference is to combine multiple methods to achieve the best possible result. This is the strategy of the work *Strategies for combining Twitter users geo-location methods* [91]. In this work the authors analyzed a bunch of the aforementioned works [26] [7] [55] [58] [94] [23] [51]. After analyzing them, they also thought that it is possible to combine location inference through friendships with other methods that exploit other types of information, such as text-based methods. Text-based methods use different sources of text information to infer users' locations. The text source can be extracted from: (i) user self-reported location field, (ii) the description profile, (iii) Twitter user name, (iv) tweets, and (v) a mix of the previous ones. Some works [90] [92] tried to improve accuracy using this combination, but did not manage to accomplish that. A work that manage to achieve better results than using the two methods individually is *Text-based Twitter user geolocation prediction* [41]. In work [91] to improve the state of the art, they use [41] and [43] text-based methods in order to actuate experimental comparisons. In [91] they evaluate the aforementioned methods by using the following metrics:

- Recall: percentage of users that the method is able to predict a location to.
- Acc@k: percentage of users whose location was inferred within a  $k$  kilometers radius of the real location.
- AUC-g: it calculates the distance between the real and predicted city groups, considering the distance from the center of the inferred city group to the center of the user real group.
- AUC-c: it calculates the distance between the center of the assigned group to the user's real city.
- Mean: arithmetic mean of the distances between the center of the groups assigned

to the user and the user’s real city.

The outcomes of the study reveal that FindMe, when coupled with the mentions network, emerges as the frontrunner in terms of accuracy-related metrics, presenting the most favorable overall results. On the other hand, the Spot method, leveraging the friendship network, showcases the highest recall performance. Nevertheless, an equilibrium between these two metrics is not strikingly achieved by either of these methods. Interestingly, the Neighborhood Vote approach, which takes the friendship network into account, attains the most noteworthy accuracy-related outcome (AUC-c of 0.58 and AUC-g of 0.81) and covers 90% of the data-set. On initial observation, Neighborhood Vote appears to hold the advantage as the optimal overall technique. However, it is essential to consider the intricacies of collecting the friendship network from social media platforms, such as Twitter, which often entails considerable resource investment. Furthermore, the efficacy of this approach may be compromised in scenarios where the network exhibits low density, thereby warranting careful consideration. In [91] the aforementioned methods are compared, leading to the following observations: they have high level of disagreement in the inferred location and they also cover very different sets of users; this opens up the opportunity to take advantage of all these methods together. So they combine the different methods using two approaches: voting committees and meta-decision trees. This is attributed to the conjecture that due to varying practices in location disclosure among users, employing a blend of diverse techniques across distinct user metadata can potentially enhance the ability to infer the location of a larger user population with heightened precision. For voting committees the authors used tree variations:

- Majority Vote (MV): it considers as the inferred location the one attributed by the highest number of base methods.
- Weighted Accuracy Vote (WAVE): similar to MV, but the contribution of different base methods to the final location is allocated by using weights for each one that are assigned according to the accuracy in the validation partition.
- Genetic Adjusted Vote (GAVE): in the same way as WAVE, it modifies the contributions of the votes of different classifiers using weights that are iteratively optimized using a genetic algorithm.

The fourth proposed strategy is grounded in Meta Decision Trees (MTDs), which share the structural framework of conventional decision trees but emphasize amalgamating outcomes from diverse classifiers. Decision trees, a form of supervised learning, are deployed for classification and regression tasks. In these trees, internal nodes encapsulate data attribute conditions, and each leaf node corresponds to a projected class for that tree



path. What distinguishes MTDs is that each leaf node denotes a method designated for forecasting the location of a new user. According to their empirical findings, GAVe and MDT strike an optimal balance between precision and recall. The result achieved using GAVe are: coverage of 98% of the totality of users without a friendship network, and accurate categorization of 61% of them in a range of 100km from their actual location. This work is relevant because it achieved nearly complete user coverage while also enhancing accuracy.

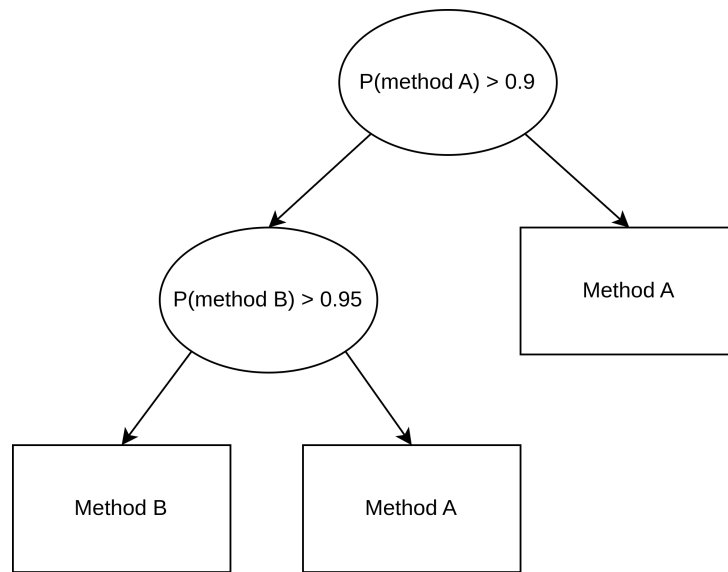


Figure 3.2: Meta Decision Tree example.

### 3.4.2. Location Inference through User Similarity

Instead of trying to localize a user from its friendship network, many work focus on finding user who share similar moving patterns and/or habits to infer location. One way to approach this is by analyzing co-location information. With co-location we refer to two users appearing in the same location at the same time (or in the same time interval). This is useful because we can extract information from users who share similar moving patterns.

The first work that this thesis presents is *Deanonymizing Mobility Traces With CoLocation Information* [53]; the primary focus of this study involves harnessing data from Twitter and Swarm to execute two distinct attacks, leveraging historical location traces for the creation of user mobility profiles. In the devised attack scenarios, the adversary functions as an observer with access to divulged mobility traces. In the first attack, the observer possesses information regarding a user’s location and co-location data within the

observed temporal span. In the second attack, the adversary augments this knowledge with access to past location traces to formulate a comprehensive user mobility profile. The adversary's primary objective revolves around identifying the trace that best aligns with the provided user information across all observed traces. For the initial attack, a maximum likelihood estimator (MLE) is employed to discern the obscured mobility trace of a user. Transitioning to the second attack, the inference problem translates into resolving a first-order Hidden Markov model (HMM). To accomplish this, the researchers employ a widely recognized inference algorithm termed iterative forward algorithm. Evaluation of the first attack demonstrates an enhancement in identification accuracy when observed locations are coupled with co-location data. In essence, a greater availability of location (co-location) information amplifies the attacker's likelihood of pinpointing the intended user's trace. Subsequent testing of the second attack reveals its effectiveness even when solely relying on co-locations. Remarkably, the attacker can successfully identify up to 17% of users within the traces using solely mobility profiles. The efficacy of the attack corresponds to the accuracy with which the constructed mobility profiles depict users' movements, underscoring the interplay between profile precision and attack success.

User's behavior habits can also be used in conjunction to other factors seen before such as textual content, social relationships and/or user similarity to predict user's location. This has been done in the work *Location Prediction in Social Networks* [62], where they combine the three factors just mentioned to predict user's current locations for tweets without any location tags. To be more specific, in this work they use 3 main factors: textual content, social relationship, and behavior habit. When examining textual content, the researchers seek out indicators and words associated with specific locations within tweets. Concerning social relationships, the focus isn't solely on online friendships, but rather on users who share a local connection determined by the similarity of their trajectories, rather than explicit friendship ties. Regarding behavioral habits, the aim is to extract patterns from historical data, such as a user's movement between places of residence. These three factors are addressed using distinct models. To address content-based analysis, a Convolutional Neural Network (CNN) is employed to establish the interrelation between textual data and geographical locations. In the social relationship model, a novel approach is introduced that quantifies the similarity between two users by evaluating the resemblance of their trajectory sequences (sequence of user's check-ins). Beyond a certain threshold, users with similar trajectories are identified, enabling the estimation of the likelihood of their presence in a specific location. As for the behavioral habit component, a Markov chain (utilizing the Monte Carlo method) is employed to predict a user's forthcoming location based on their locations. Upon evaluating experimental outcomes, it is deduced that the content-based model is most adept at handling tweets containing

location-specific terms. For the relationship-based model, varying thresholds for the similarity parameter were tested, revealing a modest enhancement in prediction accuracy with larger thresholds. In the behavior habit-based model, the dataset is divided into training and testing subsets. The accuracy of this location prediction model rises proportionally with an increase in error distance, successfully locating around 41% of predicted tweets within a 5 km radius of their actual locations. In a bid to optimize results, a linear combination model was tested, amalgamating the last two models. Notably, the fusion of the relationship-based model with the behavior habit-based model nearly doubled the prediction accuracy compared to leveraging solely the social relationship-based model. In conclusion, while the first model excels specifically with location-related word-rich tweets, a combination of the other two models can offer precise outcomes for other types of tweets. In the work *No place to hide: Inadvertent location privacy leaks on Twitter* [95] a novel inference technique called Jasoos is introduced. Jasoos employs a modified Naive Bayes approach to detect common vocabularies shared among users, considering both temporal and non-temporal aspects. Through the analysis of user sharing behavior, this method pinpoints the geographical locations of social media users.

Studying social relation affinity and moving pattern similarity can be useful not only to infer someone's location, but also to predict a user trajectory (next-location prediction). The work *Modeling user mobility for location promotion in location-based social networks* [123] introduces a distance-based mobility model to represent user check-in behavior within location-based social networks (LBSNs). To optimize location promotion, they frame it as an influence maximization problem in a LBSN. The question becomes: given a target location and an LBSN, which initial set of  $k$  users (referred to as seeds) should be advertised to effectively propagate and attract the majority of other users to visit the target location? Their evaluation involves employing Gaussian-based mobility models (GMM) and distance-based mobility models (DMM) on Gowalla and Brightkite datasets. The comprehensive experimental findings highlight that the DMM approach surpasses other existing methods by effectively capturing individual check-in behaviors. In the work *Social lstm: Human trajectory prediction in crowded spaces* [4] the authors propose a "Social long short-term memory network (LSTM)" designed for predicting human motion dynamics or future trajectories within densely populated areas. By using a "social" pooling layer which shares hidden states, Social-LSTM achieves more accurate trajectory predictions compared to current state-of-the-art methods on two distinct datasets and also demonstrates the ability in jointly predicting trajectories of pedestrians moving in groups or pairs. Another noteworthy work called *Human mobility prediction through Twitter* [22] introduces a strategy called "Similarity-based next-place Prediction from Twitter" (SimPreT) that focuses on forecasting a user's subsequent location by uti-

lizing historical data and a custom similarity function designed for trajectory patterns. This technique encompasses sequential pattern mining while incorporating the Haversine distance as the pattern similarity metric to align with the user’s ongoing trajectory. In cases where multiple patterns score the highest similarity, the approach combines various indicators encompassing user location regularities and patterns. The experimental findings reveal precision of 84%, recall of 91%, and an F1-measure of 87%, showcasing the superior performance of this method in comparison to existing state-of-the-art approaches. The authors of [33] further improved their work in the following one [19] by differentiating word importance for different locations. Due to the limited information provided by a single tweet, since users often go to the same places according to their routine or external factors, the authors recur to query expansion in order to include the user’s previous tweets as additional information. Another work is *A User Location Prediction Method Based on Similar Living Patterns* [47]; this work addresses the challenge of improving user’s location prediction accuracy, which is often hindered by the sparsity of users’ check-ins. To address this issue, the authors propose a method that initially creates vector representations of users’ living habits, allowing for clustering of users with similar living patterns. This method focuses on three concepts: check-in points, check-in traces, and users’ living patterns. These patterns encompass the most frequently visited types of Points of Interest (POIs) for each time slice, prioritizing the time a user visits a POI. To effectively capture both POI category and temporal information while measuring user similarity, the paper adopts a representation learning approach based on Global Vectors for Word Representation (GloVe) [76], embedding user life patterns into a common vector space. The derived model, referred to as POI Type to Vector (PT2V), segregates the prediction of locations into two distinct tasks: forecasting activities through POI type embedding vectors and predicting locations using POI location embedding vectors. Real user’s check-in data demonstrates that in various cases this approach consistently outperforms baseline methods.

There is a category of works that analyzes text from posts and give a central role to locations: locations are treated as pseudo-documents that encompass all tweets from users residing in those locations. To forecast a user’s home location, the pseudo-document of that user is compared to other pseudo-documents, and the locations sharing the most similar pseudo-documents are identified as prediction outcomes [121]. The authors from [108] employ a grid-based representation for locations. They construct a language model [80] for each grid using its corresponding pseudo-documents. To enhance the probability estimation for unseen words, they apply Good-Turing smoothing [36]. To quantify the likeness between location and user documents, they employ the Kullback-Leibler divergence. Inspired by [93], the following work by the same authors [109], utilizes adaptive

grids. In situations requiring the reporting of geo-coordinates rather than grids, they found that indicating the centroid of user locations within the grid offers improved accuracy compared to indicating the mid-points of the grid.

A list of all the works from this section and all of their differences is shown below in Table 3.2.

Work	Approach	Method	Data	Inferred location
[7]	Friendships	Probabilistic method	Friends' geographical distance	Home location
[55]	Friendships	Probabilistic method	Friends' geographical distance	Home location
[70]	Friendships	Probabilistic method with decision tree	Friends' geographical distance	Home location
[59]	Friendships	Probabilistic method	Friends' locations and mentioned locations	Home location
[69]	Friendships	Probabilistic method	Followers, mentions and conversations	Home location
[84]	Friendships	Probabilistic method	Friends' locations	Home location
[74]	Friendships	Bayesian network	Reported co-locations with friends	Activity location
[96]	Friendships	Dynamic Bayesian Network	Friends' location sequences	Activity location
[75]	Friendships	Bayesian hidden location inference model and a multi-factor fusion based hidden location inference model	Friends' similarity	Activity location
[58]	Friendships	Supervised version of Latent Dirichlet Allocation	Multiple home locations	Home location

Work	Approach	Method	Data	Inferred location
[10]	Friendships	Artificial neural network	Multiple explicit and implicit positions of friends	Home location
[72]	Friendships	Neural network model	Mentions network and text	Home location
[101]	Friendships	Distance Vector, Randomized Spanning Trees, Recursive Subgraph Matching heuristics	Anonymized mobility traces	Activity location
[49]	Friendships	Distance Vector, Randomized Spanning Trees, Recursive Subgraph Matching heuristics	Anonymized mobility traces	Activity location
[51]	Friendships	LP algorithm	Geometric median of location of user's neighbors	Home location
[23]	Friendships	LP algorithm	Geometric median of location of user's neighbors with weights	Home location
[113]	Friendships	LP algorithm	Geo-related content of friends	Home location
[39]	Friendships	LP algorithm	Co-locations with friends	Activity location
[104]	Friendships	LP algorithm and representation learning	Mentions in user-generated texts	Home location
[86]	Friendships	LP algorithm and text-based methods	Mentions network and text	Home location
[85]	Friendships	LP algorithm	Mentions network	Home location

Work	Approach	Method	Data	Inferred location
[26]	Friendships	Majority Vote	Friends' locations	Activity location
[34]	Friendships	LocusRank algorithm	Friends' self reported locations	Home location
[94]	Friendships	SVM	Number of friends per city	Home location
[11]	Friendships	Social Tie Factor Graph Model	Following network from Twitter	Home location
[38]	Friendships	Influence model and global iteration algorithm	Number of common friends	Home location
[114]	Friendships	Landmark Mixture Model	Friends that live close to the user	Home location
[91]	Friendships	Mix of more methods using voting committees and decision trees	Output of other methods	Home location
[53]	User similarity	Maximum likelihood estimator and Hidden Markov model	Co-locations and published mobility traces	Activity location
[62]	User similarity	Convolutional Neural Network, Markov chain based on Monte Carlo	Textual content, behavior habits, similarity between user's trajectories	Activity location
[95]	User similarity	Naive Bayes	Shared vocabularies between users	Activity location
[123]	User similarity	Gaussian-based mobility model	user's check-ins from LBSN	Next location
[4]	User similarity	Social long short-term memory network	Individual or group trajectories	Next location
[22]	User similarity	Similarity function	Trajectory patterns	Next location

<b>Work</b>	<b>Approach</b>	<b>Method</b>	<b>Data</b>	<b>Inferred location</b>
[47]	User similarity	Representation learning	user's check-ins	Next location
[19]	User similarity	Collaborative filtering	Mobility behavior of users	Activity location
[108]	User similarity	Pseudo-documents	Tweets from users	Home location
[109]	User similarity	Pseudo-documents	Tweets from users	Home location

Table 3.2: Works regarding Localization through Social Network



# 4 | Friendship Inference through Location Data

In this thesis the problem of inferring social relationships by exploiting location data from social media is introduced. This problem is relevant to location inference because it tries to solve the exact opposite problem: in location inference we can exploit friends' data to infer location more accurately, in friendship inference we can exploit location data to discover users' relationships. Friendships inherently reflect social interactions and patterns, which have a significant impact on individuals' movement and activities. When users' friendships are taken into account, it can enhance the precision of location predictions by leveraging shared patterns and locations among friends. In fact, addressing friendship inference through location data enriches the understanding of how social relationships can contribute to enhancing the accuracy and efficacy of location predictions. This section is divided in three parts: section 4.1 covers methods that use check-in data to infer friendships; section 4.2 covers methods that exploit trajectories to infer friendships; section 4.3 covers methods that use multiple sources including location data.

## 4.1. Friendship Inference through Check-Ins

Approaches grounded in check-in location data emphasize individual check-in points rather than sequences or trajectories. These methods discern users' inclinations or patterns in the real world through their recorded check-in locations or areas [107].

In study [46], it was observed that around 30% of new connections emerged among users who shared common places they visited. As a result, several investigations evaluated users' proximity by extracting co-location attributes, encompassing factors like co-location frequency, proximity of significant locations, and the likelihood of co-location. In another study [105], the researcher introduces a range of check-in location attributes, including GeoDist and check-in observations, and then conducts a comparative analysis of their predictive efficacy. To harness diverse information sources, multiple researches effort to

amalgamate online and offline attributes, subsequently evaluating user similarities through the definition of multiple features. In a separate study [25], a trio of network types were established: social network, co-location network, and co-located friend networks. Within each network, approximately 67 attributes were defined across five categories to delineate distinct characteristics. Furthermore, numerous user-generated contents, such as posts related to Points of Interest (POIs), also proliferated in Location-Based Social Networks (LBSNs). To predict social relations, the paper titled *vec2Link: Unifying heterogeneous data for social link prediction* [122] integrated offline check-in behavior and users' online behavior. The *vec2link* framework employs a neural network for embedding user social relation and uses a location sensitive hash for efficient convolutional network learning. The study *Inferring social ties from geographic coincidences* [24], employs geographic coordinates of users in a social network to infer social relationships via a co-occurrence model. *Mobility intention-based relationship inference from spatiotemporal data* [115] introduced the MIRI model, inferring friend relationships through mobility intention dyads as features. This model leveraged spatiotemporal data to enhance relationship inference. The scope of this paper is to solve the following problem: in spatiotemporal datasets generated by IOT devices in public locations, it's difficult to distinguish co-locations between friends and strangers. The work *Exploiting place features in link prediction on location-based social networks* [98] created a social graph based on users who visited the same locations, forming a supervised learning model. Their findings indicated that approximately 30% of new social links are established among users who frequent the same places, termed as "place-friends" in the study. The primary objective of this study is to uncover potential future connections among users.

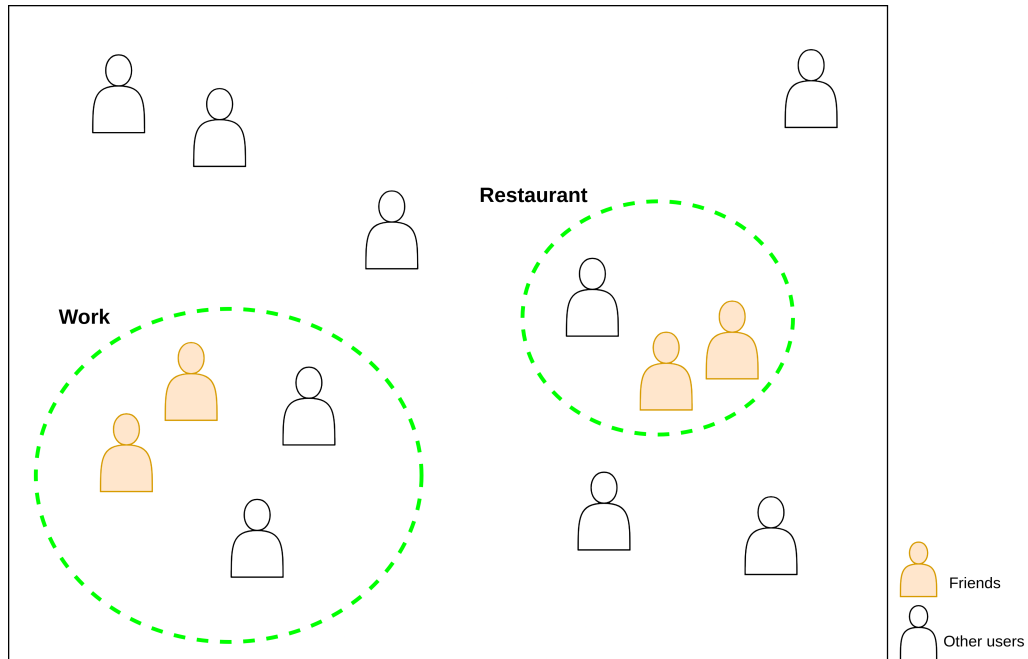


Figure 4.1: Co-locations between friends example.

## 4.2. Friendship Inference through Trajectories

While the incorporation of check-in location data indeed enhances the precision of friendship prediction, certain constraints remain evident. Firstly, the method may yield inaccurate estimations of actual social strength for friends who infrequently share locations or for unrelated individuals who frequently coincide at the same places. Secondly, these investigations overlook the dynamics of users' overall trajectory patterns as they shift in response to different lifestyles. Hence, in contemporary times, trajectory similarity has been harnessed for the scrutiny of user mobility homophily. Associated techniques delve into the underlying trajectory patterns within user trajectories, transcending mere geositions, to explore and evaluate the likeness in mobility among users. This notion is rooted in the belief that trajectory sequences mirror users' distinctive "lifestyles", thereby fostering connections between individuals who share similar "lifestyles" [107].

In the study outlined in [116], the researchers utilizes the tags corresponding to landmarks that users traversed as semantic labels for their trajectories. This approach resulted in a trajectory representation such as "school  $\rightarrow$  park  $\rightarrow$  restaurant", forming a semantic sequence. Subsequently, they extracted the most significant semantic trajectory pattern from each user's trajectory, employing it to quantify the likeness between trajectories. In a separate study, detailed in [112], the authors transformed each stay region into a distinct feature vector and organized these stay points into distinct categories, each endowed with

a distinct semantic implication. These categories collectively formed a semantic location history (SLH), wherein user trajectories were encoded. Finally, SLH permits to measure user similarity. As a supplementary method to trajectory similarity based methods that can be used to improve friendship inference is topic-based methods. This because they can extract location semantics and users' preferences. In [119], the researchers introduces a probabilistic generative model, which facilitates the discovery of patterns determined by lifestyle within users' trajectories. This model encompasses a series of features: (i) users' preferences, (ii) interdependence between diverse locations, (iii) service duration, (iv) users' lifestyle. The research presented in *Inferring online social ties from offline geographical activities* [45] introduced the O2O-Inf framework, which capitalizes on users' offline geographical activities such as check-in records and meetings to infer online social relationships. The researchers devised a linkage graph to illustrate relationships between nodes and utilized a graph-based SSL method. This iterative process computes the probability of nodes becoming friends by taking into account adjacent nodes within the linkage graph. In *walk2friends: Inferring social links from mobility profiles* [6], the authors transferred a method based on deep learning for inferring user's location, to the inference of social relationships. Random walk traces in users' location graph represents their neighbors, and similarity measurements between users help in predicting social connections. *Graph convolutional networks on user mobility heterogeneous graphs for social relationship inference* [111] harnesses a graph convolutional network (GCN) to construct a heterogeneous mobility graph with an unsupervised method that analyzes human trajectories. In order to infer social relationships, this graph covers a series of features: (i) user-user meeting, (ii) social graph, (iii) user-location bipartite, (iv) location-location co-occurrence graphs. *CIFEF: Combining implicit and explicit features for friendship inference in location-based social networks* [42] focuses on learning track features from check-in sequences, proposing the CIFEF method. This approach combines implicit features from weekdays' and weekends' trajectory patterns with a new explicit feature based on shared locations, effectively inferring friendships even between users with minimal co-occurrence.

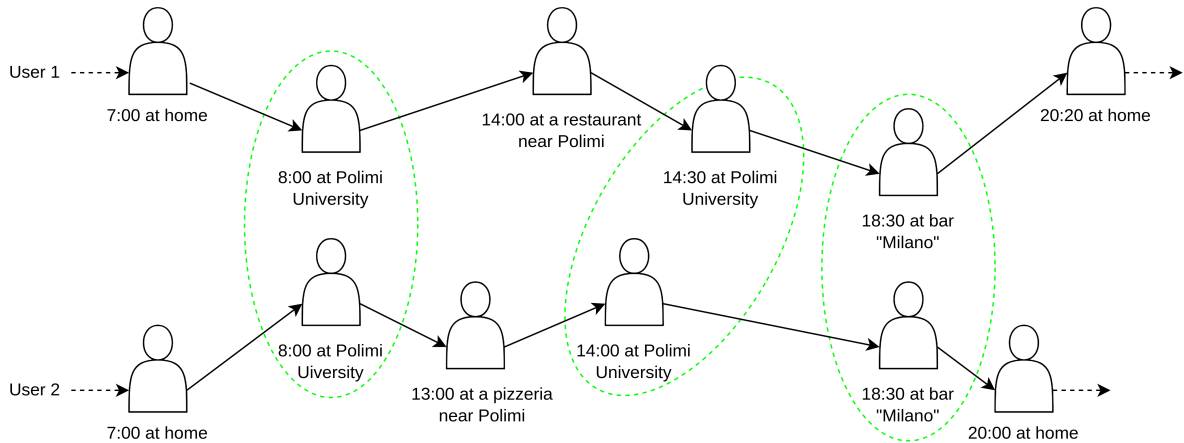


Figure 4.2: Trajectory example of two students that are friends.

### 4.3. Friendship Inference through Multiple Sources Including Location Data

As exemplified many times in the location inference section, location information can be used in conjunction with other information sources such as published images, tags etc. to better infer social relationships. This mix of different sources and types of information can often yield to better results.

The authors from [120] introduced a two-stage deep learning framework, known as TDFI, designed for friendship inference. This framework employs an extended adjacency matrix (EAM) to capture diverse information from multiple sources. Utilizing the enhanced deep Siamese network IDSN, the fusion features' similarity is measured to determine users' friendships. In another study, [63] the authors proposed an inference model that combines random forest, SVM, and naive Bayes techniques to infer friendships. Relating to location-based social networks, this fusion of multi-source data encompasses different features: (i) user social topology, (ii) location categories, (iii) check-in locations. The work [89] employees long-term multimodal information to deduce social relationships. This approach integrates various sources, including images, tags, titles, geographic locations, and established friendships, to enhance the accuracy of friendship inference. *EBM: An entropy-based model to infer social strength from spatiotemporal data* [78] proposed an entropy-based model using spatial-temporal data to infer social connections based on diversity and weighted frequency. Other than inferring social connections, this paper also infers the strength of social relationships by analyzing people's co-locations in space and time. This work highlights the efficacy of merging location-based data with other data types such as temporal data and social attributes to deduce social connections.



# 5 | Discussion and Open Research Questions

This chapter will provide observations about the current state of the art, what its current limitations are, and which research questions remain open. The observations made are divided in two sections to distinguish between the state of the art regarding location privacy protection and regarding location inference.

Regarding location privacy, this thesis analyzed the two privacy models:  $k$ -anonymity and geo-indistinguishability. These two models are the main guidelines that social media and location services providers follow when trying to protect users privacy, but they are not bullet-proof and each of them has their own limitations.  $K$ -anonymity does not defend against background knowledge attacks. If attackers have access to external information or auxiliary data sources, they may be able to de-anonymize individuals in the dataset by matching it with external data, potentially revealing sensitive information. For this reason, mixing different types of data to infer locations is very effective against  $k$ -anonymity. When it comes to geo-indistinguishability, instead, its main limitation is the excessive decrease in data utility when applied. Typically, to achieve geo-indistinguishability, the common approach is to use perturbation-based mechanisms, which introduce noise, this results in a utility loss when it comes to the data in question. Having made these two observations, an open research direction is to invest more efforts and researches into hiding correlations between location data and temporal data; the reason behind this observation is that a lot of information is derived from spatio-temporal data and not only from location data. For example, it is possible to tell whether a bar is visited during work breaks or during evenings by the temporal aspect, and this can be used to infer information about users such as workplaces and so on. This research direction should improve the current state of the art in two ways: hiding does not decrease utility of the data as aggressively as perturbation-based techniques do, especially when hiding links between location data and temporal data; furthermore, this would also overcome the previously mentioned limitation of  $k$ -anonymity, since we are removing additional information from location data. Additionally, putting more focus on to contextual anonymization is an open question in

the location privacy field: developing methods to anonymize location data (for example by suppressing or sampling) in conjunction with contextual information would provide useful but privacy-preserving insights. Additionally, data can be aggregated with contextual attributes: that is to group location data based on the identified contextual attributes; for example, by aggregating location data by time intervals or venue categories. Improving data utility while preserving location privacy remains an open research question. Improving utility can be achieved by utilizing local differential privacy: it consists in adding noise to individual data points before aggregation. This can enhance privacy without sacrificing too much utility. Another observation can be made from the following fact: one of the most common approaches to infer location using social relationships is label propagation. Since the Label Propagation algorithm focuses on analyzing graph data, LPPMs focused on graph data and their effectiveness against the LP algorithm should be further investigated, to counter this type of location inference. There is place for more investigation on this specific subject. Techniques such as random edge deletion (hiding), random edge addition and random edge switching need to be tweaked to increase privacy protection against LP algorithm attacks. Since the LP algorithm is non deterministic, increasing privacy protection on graph data should make the optimization of the LP algorithm more difficult. This because there would be more cases where it is not certain which label an entity in the graph should receive. The degree of privacy protection achieved against LP algorithm attacks in relation to the number of altered edges and the loss of data utility is an open question for research.

Moving on to the subject of location inference attacks, from the research conducted in this thesis, a common trend that can be noticed is the fact the researchers mention only the social media that the data came from, but don't make any comments about which LPPMs were implemented in that specific data-set. The most common observation that researchers do when inferring location is the sparsity of geo-tagged posts, or the noise in self-reported locations. This gives an interesting perspective, because it shows the scenario where an attacker has no knowledge regarding the protection mechanisms used. That said, it would benefit the research in the state of the art if more focus is put on the correlations between the LPPMs and the location inference techniques, from the perspective of location inference researchers. There are too few papers that analyze this subject in the works specifically devoted to location inference, at the moment. Also another currently open research question that can benefit from further research is the use of deep learning methods (artificial neural networks etc.) This research direction has been already explored by various works, but it is very likely that it will be pursued more and more. Deep learning should be more and more dominant in location inference approaches for



its capacity of extracting features from raw data. Moreover, deep learning models inherently learn meaningful representations of data, which can be valuable for understanding underlying patterns and relationships. Another reason is their scalability: Deep learning models can scale with the availability of more data and computational resources. This scalability allows them to handle large datasets and complex problems effectively. Particular relevance has the fact that deep learning models can handle large datasets: with the continuous rising of social media usage, datasets will only increase in dimensions. In particular, graph neural networks are particularly useful for location inference when considering the friendship network, because they can capture the complex relationships between users through graphs. Another open research question is the use of data from multiple social media to better improve location inference. Using data from different social media will solve in part the problem of sparse data (one of the major obstacles in location inference): by including different sources, the number of information can be significantly increased, with the cost of a significant increase in computational complexity. The additional data sources would make the rise in computational complexity worth it; this could permit to study how users' behavior and relationships differ across various social media platforms. Developing techniques that can generalize across platforms would be valuable. Another open research direction is to take inspiration from the approach presented in [91], where different methods were combined using voting systems and meta decision trees. This approach can be further extended by adding different ways of combining and choosing results from different methods; for example, averaging results can be used in cases where different methods achieve different results for the same user and, for cases where methods cover different sets of data (or some are better than others for certain sets of users), other types of decision support systems and/or ensemble methods can be further investigated. Another interesting open research direction is mobile and IoT integration: that is to incorporate data from mobile devices and the Internet of Things (IoT) into location inference models to create a more comprehensive view of users' locations. This integration will benefit location inference by providing more location data from mobile devices in real time. This could also be used for real time location inference; this could be useful for emergency response or location-based marketing. Incorporating multi-modal data remain an open research question: exploring how other types of data, such as images, audio, or video shared within social networks, can enhance location inference. Integrating multi-modal data could provide richer context for understanding users' behavior. In particular, the effectiveness of computer vision technologies to extract location information from images and videos shared on social media platforms is an interesting research direction. Lastly, using contextualized and advanced language models such as GPT-3 or similar could help to extract location-related information from textual data. Since these

language models are very recent and still under improvement, their potential for location inference remains an open research question.

## 6 | Conclusions

This work focused on various subjects regarding location privacy. While the scope of this thesis was to give a complete overview about the state of the art regarding location inference through friendships, in order to give a more complete and broader description of the subject, this survey also covered the various aspect related to location data and location privacy. To achieve this, it also analyzed the topic of location privacy protection mechanisms, location inference in general, and friendship inference using location data from social media. In Chapter 1, this work covered works that focused on the same or similar subjects covered in this thesis. This by describing their focuses, achievements, and differences from this survey. Then, in Chapter 2, this work covered the topic about location privacy preserving mechanisms, highlighting their characteristics, how they work, how to evaluate them, and how effective they are. This chapter preceded the chapter about location inference in order to give the reader a general understanding of the protection mechanisms that the location inference attacks analyzed in the following chapters have to deal with. Chapter 3 analyzed the topic of location inference, dividing the examined works in three categories based on the type of location inferred: home location, activity location and next location. Then, Chapter 3 covered the topic of inferring location through users' social network (users' relationships such as friends, and people surrounding them), by differentiating methods that use data from friends from methods that exploit similarity between users' behaviors. In Chapter 4, the subject of inferring friendships using location data was briefly covered. This was to highlight the opposite side of the problem covered in Chapter 3 (location inference through friends' social media). Then, Chapter 5 presented observations and open research questions.



## Bibliography

- [1] *Inferring the Origin Locations of Tweets with Quantitative Confidence*, CSCW '14, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450325400. doi: 10.1145/2531602.2531607. URL <https://doi.org/10.1145/2531602.2531607>.
- [2] R. Ahuja, G. Ghinita, N. Krishna, and C. Shahabi. Protecting against inference attacks on co-location data. In *2019 28th International Conference on Computer Communication and Networks (ICCCN)*, pages 1–11. IEEE, 2019.
- [3] N. Al Hasan Haldar, J. Li, M. Reynolds, T. Sellis, and J. X. Yu. Location prediction in large-scale social networks: an in-depth benchmarking study. *The VLDB Journal*, 28:623–648, 2019.
- [4] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016.
- [5] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Ge-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 901–914, 2013.
- [6] M. Backes, M. Humbert, J. Pang, and Y. Zhang. walk2friends: Inferring social links from mobility profiles. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1943–1957, 2017.
- [7] L. Backstrom, E. Sun, and C. Marlow. Find me if you can: improving geographical prediction with social and spatial proximity. In *Proceedings of the 19th international conference on World wide web*, pages 61–70, 2010.
- [8] G. Beigi and H. Liu. Privacy in social media: Identification, mitigation and applications. *arXiv preprint arXiv:1808.02191*, 2018.

- [9] G. Beigi and H. Liu. A survey on privacy in social media: Identification, mitigation, and applications. *ACM Transactions on Data Science*, 1(1):1–38, 2020.
- [10] W. Chanthaweethip, X. Han, N. Crespi, Y. Chen, R. Farahbakhsh, and A. Cuevas. “current city” prediction for coarse location based applications on facebook. In *2013 IEEE Global Communications Conference (GLOBECOM)*, pages 3188–3193. IEEE, 2013.
- [11] J. Chen, Y. Liu, and M. Zou. From tie strength to function: Home location estimation in social network. In *2014 IEEE Computers, Communications and IT Applications Conference*, pages 67–71. IEEE, 2014.
- [12] M. Chen, X. Yu, and Y. Liu. Mpe: a mobility pattern embedding model for predicting next locations. *World Wide Web*, 22(6):2901–2920, Nov 2019. ISSN 1573-1413. doi: 10.1007/s11280-018-0616-8. URL <https://doi.org/10.1007/s11280-018-0616-8>.
- [13] M. Chen, Y. Zuo, X. Jia, Y. Liu, X. Yu, and K. Zheng. Cem: A convolutional embedding model for predicting next locations. *IEEE Transactions on Intelligent Transportation Systems*, 22(6):3349–3358, 2021. doi: 10.1109/TITS.2020.2983647.
- [14] Y. Chen, J. Zhao, X. Hu, X. Zhang, Z. Li, and T.-S. Chua. From interest to function: Location estimation in social media. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 27, pages 180–186, 2013.
- [15] Z. Cheng, J. Caverlee, and K. Lee. You are where you tweet: a content-based approach to geo-locating twitter users. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 759–768, 2010.
- [16] Z. Cheng, J. Caverlee, H. Barthwal, and V. Bachani. Who is the barbecue king of texas? a geo-spatial approach to finding local experts on twitter. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 335–344, 2014.
- [17] S.-B. Cho. Exploiting machine learning techniques for location recognition and prediction with smartphone logs. *Neurocomputing*, 176:98–106, 2016. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2015.02.079>. URL <https://www.sciencedirect.com/science/article/pii/S092523121500569X>. Recent Advancements in Hybrid Artificial Intelligence Systems and its Application to Real-World Problems.
- [18] W.-H. Chong and E.-P. Lim. Exploiting contextual information for fine-grained

- tweet geolocation. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 11, pages 488–491, 2017.
- [19] W.-H. Chong and E.-P. Lim. Tweet geolocation: Leveraging location, user and peer signals. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 1279–1288, 2017.
- [20] C.-Y. Chow, M. F. Mokbel, and W. G. Aref. Casper\* query processing for location services without compromising privacy. *ACM Transactions on Database Systems (TODS)*, 34(4):1–48, 2009.
- [21] C. Comito. Where are you going? next place prediction from twitter. In *2017 IEEE international conference on data science and advanced analytics (DSAA)*, pages 696–705. IEEE, 2017.
- [22] C. Comito. Human mobility prediction through twitter. *Procedia computer science*, 134:129–136, 2018.
- [23] R. Compton, D. Jurgens, and D. Allen. Geotagging one hundred million twitter accounts with total variation minimization. In *2014 IEEE international conference on Big data (big data)*, pages 393–401. IEEE, 2014.
- [24] D. J. Crandall, L. Backstrom, D. Cosley, S. Suri, D. Huttenlocher, and J. Kleinberg. Inferring social ties from geographic coincidences. *Proceedings of the National Academy of Sciences*, 107(52):22436–22441, 2010.
- [25] J. Cranshaw, E. Toch, J. Hong, A. Kittur, and N. Sadeh. Bridging the gap between physical location and online social networks. In *Proceedings of the 12th ACM international conference on Ubiquitous computing*, pages 119–128, 2010.
- [26] C. A. Davis Jr, G. L. Pappa, D. R. R. De Oliveira, and F. de L. Arcanjo. Inferring the location of twitter messages based on user relationships. *Transactions in GIS*, 15(6):735–751, 2011.
- [27] M. Dredze, M. Osborne, and P. Kambadur. Geolocation for twitter: Timing matters. In *Proceedings of the 2016 conference of the north American chapter of the association for computational linguistics: human language technologies*, pages 1064–1069, 2016.
- [28] C. Dwork. Differential privacy. In *Automata, Languages and Programming: 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II 33*, pages 1–12. Springer, 2006.

- [29] H. Efstathiades, D. Antoniadou, G. Pallis, and M. D. Dikaiakos. Identification of key locations based on online social network activity. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015*, pages 218–225, 2015.
- [30] J. Eisenstein, B. O’Connor, N. A. Smith, and E. Xing. A latent variable model for geographic lexical variation. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 1277–1287, 2010.
- [31] J. Eisenstein, A. Ahmed, and E. P. Xing. Sparse additive generative models of text. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 1041–1048, 2011.
- [32] W. Enck, P. Gilbert, S. Han, V. Tendulkar, B.-G. Chun, L. P. Cox, J. Jung, P. McDaniel, and A. N. Sheth. Taintdroid: an information-flow tracking system for real-time privacy monitoring on smartphones. *ACM Transactions on Computer Systems (TOCS)*, 32(2):1–29, 2014.
- [33] D. Flatow, M. Naaman, K. E. Xie, Y. Volkovich, and Y. Kanza. On the accuracy of hyper-local geotagging of social media content. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM ’15*, page 127–136, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450333177. doi: 10.1145/2684822.2685296. URL <https://doi.org/10.1145/2684822.2685296>.
- [34] M. Ghufraan, G. Quercini, and N. Bennacer. Toponym disambiguation in online social network profiles. In *Proceedings of the 23rd SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pages 1–10, 2015.
- [35] A. Gobezie and M. Fufa. A survey on next location prediction techniques, applications, and challenges. *EURASIP Journal on Wireless Communications and Networking*, 2022, 03 2022. doi: 10.1186/s13638-022-02114-6.
- [36] I. J. Good. The population frequencies of species and the estimation of population parameters. *Biometrika*, 40(3-4):237–264, 1953.
- [37] M. Gruteser and D. Grunwald. Anonymous usage of location-based services through spatial and temporal cloaking. In *Proceedings of the 1st international conference on Mobile systems, applications and services*, pages 31–42, 2003.
- [38] Y. Gu, J. Song, W. Liu, and L. Zou. Hlgps: a home location global positioning sys-



- tem in location-based social networks. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 901–906. IEEE, 2016.
- [39] N. A. H. Haldar, M. Reynolds, Q. Shao, C. Paris, J. Li, and Y. Chen. Activity location inference of users based on social relationship. *World Wide Web*, 24(4): 1165–1183, 2021.
- [40] B. Han, P. Cook, and T. Baldwin. A stacking-based approach to twitter user geolocation prediction. In *Proceedings of the 51st annual meeting of the association for computational linguistics: system demonstrations*, pages 7–12, 2013.
- [41] B. Han, P. Cook, and T. Baldwin. Text-based twitter user geolocation prediction. *Journal of Artificial Intelligence Research*, 49:451–500, 2014.
- [42] C. He, C. Peng, N. Li, X. Chen, Z. Yang, and Z. Hu. Cifef: Combining implicit and explicit features for friendship inference in location-based social networks. In *Knowledge Science, Engineering and Management: 13th International Conference, KSEM 2020, Hangzhou, China, August 28–30, 2020, Proceedings, Part II 13*, pages 168–180. Springer, 2020.
- [43] B. Hecht, L. Hong, B. Suh, and E. H. Chi. Tweets from justin beiber’s heart: the dynamics of the location field in user profiles. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 237–246, 2011.
- [44] L. Hong, A. Ahmed, S. Gurumurthy, A. J. Smola, and K. Tsioutsoulouklis. Discovering geographical topics in the twitter stream. In *Proceedings of the 21st international conference on World Wide Web*, pages 769–778, 2012.
- [45] H.-P. Hsieh and C.-T. Li. Inferring online social ties from offline geographical activities. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2): 1–21, 2019.
- [46] H.-P. Hsieh, R. Yan, and C.-T. Li. Where you go reveals who you know: Analyzing social ties from millions of footprints. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pages 1839–1842, 2015.
- [47] J. Hu, H. Chen, and Q. Xiao. A user location prediction method based on similar living patterns. In *Fuzzy Systems and Data Mining VIII*, pages 41–49. IOS Press, 2022.
- [48] S. Inc. The heartbleed bug, 2017. URL <http://heartbleed.com/>.
- [49] S. Ji, W. Li, M. Srivatsa, J. S. He, and R. Beyah. General graph data de-

- anonymization: From mobility traces to social networks. *ACM Transactions on Information and System Security (TISSEC)*, 18(4):1–29, 2016.
- [50] H. Jiang, J. Li, P. Zhao, F. Zeng, Z. Xiao, and A. Iyengar. Location privacy-preserving mechanisms in location-based services: A comprehensive survey. *ACM Computing Surveys (CSUR)*, 54(1):1–36, 2021.
- [51] D. Jurgens. That’s what friends are for: Inferring location in online social media platforms based on social relationships. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 7, pages 273–282, 2013.
- [52] D. Jurgens, T. Finethy, J. McCorriston, Y. Xu, and D. Ruths. Geolocation prediction in twitter using social networks: A critical analysis and review of current practice. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 9, pages 188–197, 2015.
- [53] Y. Khazbak and G. Cao. Deanonimizing mobility traces with co-location information. In *2017 IEEE Conference on Communications and Network Security (CNS)*, pages 1–9. IEEE, 2017.
- [54] H. Kido, Y. Yanagisawa, and T. Satoh. Protection of location privacy using dummies for location-based services. In *21st International conference on data engineering workshops (ICDEW’05)*, pages 1248–1248. IEEE, 2005.
- [55] L. Kong, Z. Liu, and Y. Huang. Spot: Locating social media users based on social network context. *Proceedings of the VLDB Endowment*, 7(13):1681–1684, 2014.
- [56] J. Leyden. Dark net linkedin sale looks like the real deal, 2016. URL <https://www.theregister.co.uk/2016/05/18/linkedin/>.
- [57] N. Li, T. Li, and S. Venkatasubramanian. t-closeness: Privacy beyond k-anonymity and l-diversity. In *2007 IEEE 23rd international conference on data engineering*, pages 106–115. IEEE, 2006.
- [58] R. Li, S. Wang, and K. C.-C. Chang. Multiple location profiling for users and relationships from social network and content. *arXiv preprint arXiv:1208.0288*, 2012.
- [59] R. Li, S. Wang, H. Deng, R. Wang, and K. C.-C. Chang. Towards social user profiling: unified and discriminative influence model for inferring home locations. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1023–1031, 2012.

- [60] W. Li, P. Serdyukov, A. P. de Vries, C. Eickhoff, and M. Larson. The where in the tweet. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 2473–2476, 2011.
- [61] Q. Liu, S. Wu, L. Wang, and T. Tan. Predicting the next location: A recurrent model with spatial and temporal contexts. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1), Feb. 2016. doi: 10.1609/aaai.v30i1.9971. URL <https://ojs.aaai.org/index.php/AAAI/article/view/9971>.
- [62] R. Liu, G. Cong, B. Zheng, K. Zheng, and H. Su. Location prediction in social networks. In *Web and Big Data: Second International Joint Conference, APWeb-WAIM 2018, Macau, China, July 23-25, 2018, Proceedings, Part II 2*, pages 151–165. Springer, 2018.
- [63] H. Luo, B. Guo, Z. Wang, Y. Feng, et al. Friendship prediction based on the fusion of topology and geographical features in lbsn. In *2013 IEEE 10th International Conference on High Performance Computing and Communications & 2013 IEEE International Conference on Embedded and Ubiquitous Computing*, pages 2224–2230. IEEE, 2013.
- [64] X. Luo, Y. Qiao, C. Li, J. Ma, and Y. Liu. An overview of microblog user geolocation methods. *Information processing & management*, 57(6):102375, 2020.
- [65] J. Lv, Q. Li, Q. Sun, and X. Wang. T-conv: A convolutional neural network for multi-scale taxi trajectory prediction. In *2018 IEEE International Conference on Big Data and Smart Computing (BigComp)*, pages 82–89, 2018. doi: 10.1109/BigComp.2018.00021.
- [66] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian. l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):3–es, 2007.
- [67] J. Mahmud, J. Nichols, and C. Drews. Where is this tweet from? inferring home locations of twitter users. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 6, pages 511–514, 2012.
- [68] J. Mahmud, J. Nichols, and C. Drews. Home location identification of twitter users. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 5(3):1–21, 2014.
- [69] J. McGee, J. A. Caverlee, and Z. Cheng. A geographic study of tie strength in social media. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 2333–2336, 2011.

- [70] J. McGee, J. Caverlee, and Z. Cheng. Location prediction in social media based on tie strength. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pages 459–468, 2013.
- [71] Y. Michalevsky, A. Schulman, G. A. Veerapandian, D. Boneh, and G. Nakibly. Powerspy: Location tracking using mobile device power analysis. In *24th {USENIX} Security Symposium ({USENIX} Security 15)*, pages 785–800, 2015.
- [72] Y. Miura, M. Taniguchi, T. Taniguchi, and T. Ohkuma. Unifying text, metadata, and user network representations with a neural network for geolocation prediction. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1260–1272, 2017.
- [73] S. Narain, T. D. Vo-Huu, K. Block, and G. Noubir. Inferring user routes and locations using zero-permission mobile sensors. In *2016 IEEE Symposium on Security and Privacy (SP)*, pages 397–413. IEEE, 2016.
- [74] A.-M. Olteanu, K. Huguenin, R. Shokri, M. Humbert, and J.-P. Hubaux. Quantifying interdependent privacy risks with location data. *IEEE Transactions on Mobile Computing*, 16(3):829–842, 2016.
- [75] X. Pan, W. Chen, and L. Wu. Mobile user location inference attacks fusing with multiple background knowledge in location-based social networks. *Mathematics*, 8(2):262, 2020.
- [76] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [77] S. Perez. Recently confirmed myspace hack could be the largest yet, 2016. URL <https://techcrunch.com/2016/05/31/recently-confirmed-myspace-hack-could-be-the-largest-yet/>.
- [78] H. Pham, C. Shahabi, and Y. Liu. Ebm: an entropy-based model to infer social strength from spatiotemporal data. In *Proceedings of the 2013 ACM SIGMOD international conference on management of data*, pages 265–276, 2013.
- [79] Y. Piao, K. Ye, and X. Cui. Privacy inference attack against users in online social networks: A literature review. *IEEE Access*, 9:40417–40431, 2021. doi: 10.1109/ACCESS.2021.3064208.
- [80] J. M. Ponte and W. B. Croft. A language modeling approach to information re-

- trieval. In *ACM SIGIR Forum*, volume 51, pages 202–208. ACM New York, NY, USA, 2017.
- [81] A. Poulston, M. Stevenson, and K. Bontcheva. Hyperlocal home location identification of twitter profiles. In *Proceedings of the 28th ACM conference on hypertext and social media*, pages 45–54, 2017.
- [82] V. Primault, A. Boutet, S. B. Mokhtar, and L. Brunie. The long road to computational location privacy: A survey. *IEEE Communications Surveys & Tutorials*, 21(3):2772–2793, 2018.
- [83] A. Pyrgelis, C. Troncoso, and E. De Cristofaro. Measuring membership privacy on aggregate location time-series. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 4(2):1–28, 2020.
- [84] Y. Qian, J. Tang, Z. Yang, B. Huang, W. Wei, and K. M. Carley. A probabilistic framework for location inference from social media. *arXiv preprint arXiv:1702.07281*, 2017.
- [85] A. Rahimi, T. Cohn, and T. Baldwin. Twitter user geolocation using a unified text and network prediction model. *arXiv preprint arXiv:1506.08259*, 2015.
- [86] A. Rahimi, D. Vu, T. Cohn, and T. Baldwin. Exploiting text and network context for geolocation of social media users. *arXiv preprint arXiv:1506.04803*, 2015.
- [87] A. Rahimi, T. Baldwin, and T. Cohn. Continuous representation of location for geolocation and lexical dialectology using mixture density networks. *arXiv preprint arXiv:1708.04358*, 2017.
- [88] A. Rahimi, T. Cohn, and T. Baldwin. A neural model for user geolocation and lexical dialectology. *arXiv preprint arXiv:1704.04008*, 2017.
- [89] T. Rahman, M. Fritz, M. Backes, and Y. Zhang. Everything about you: A multimodal approach towards friendship inference in online social networks. *arXiv preprint arXiv:2003.00996*, 2020.
- [90] K. Ren, S. Zhang, and H. Lin. Where are you settling down: Geo-locating twitter users based on tweets and social networks. In *Information Retrieval Technology: 8th Asia Information Retrieval Societies Conference, AIRS 2012, Tianjin, China, December 17-19, 2012. Proceedings 8*, pages 150–161. Springer, 2012.
- [91] S. Ribeiro and G. L. Pappa. Strategies for combining twitter users geo-location methods. *GeoInformatica*, 22:563–587, 2018.

- [92] E. Rodrigues, R. Assunção, G. L. Pappa, D. Renno, and W. Meira Jr. Exploring multiple evidence to infer users' location in twitter. *Neurocomputing*, 171:30–38, 2016.
- [93] S. Roller, M. Speriosu, S. Rallapalli, B. Wing, and J. Baldrige. Supervised text-based geolocation using language models on an adaptive grid. In *Proceedings of the 2012 joint conference on empirical methods in natural language processing and computational natural language learning*, pages 1500–1510, 2012.
- [94] D. Rout, K. Bontcheva, D. Preotjiuc-Pietro, and T. Cohn. Where's@ wally? a classification approach to geolocating users based on their social ties. In *Proceedings of the 24th ACM Conference on Hypertext and Social Media*, pages 11–20, 2013.
- [95] J. Rusert, O. Khalid, D. Hong, Z. Shafiq, and P. Srinivasan. No place to hide: Inadvertent location privacy leaks on twitter. *Proc. Priv. Enhancing Technol.*, 2019 (4):172–189, 2019.
- [96] A. Sadilek, H. Kautz, and J. P. Bigham. Finding your friends and following them to where you are. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 723–732, 2012.
- [97] A. Sassi, M. Brahim, W. Bechkit, and A. Bachir. Location embedding and deep convolutional neural networks for next location prediction. In *2019 IEEE 44th LCN Symposium on Emerging Topics in Networking (LCN Symposium)*, pages 149–157, 2019. doi: 10.1109/LCNSymposium47956.2019.9000680.
- [98] S. Scellato, A. Noulas, and C. Mascolo. Exploiting place features in link prediction on location-based social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1046–1054, 2011.
- [99] A. Schulz, A. Hadjakos, H. Paulheim, J. Nachtwey, and M. Mühlhäuser. A multi-indicator approach for geolocalization of tweets. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 7, pages 573–582, 2013.
- [100] A. R. Shahid, N. Pissinou, S. Iyengar, and K. Makki. Check-ins and photos: Spatiotemporal correlation-based location inference attack and defense in location-based social networks. In *2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, pages 1852–1857. IEEE, 2018.

- [101] M. Srivatsa and M. Hicks. Deanononymizing mobility traces: Using social network as a side-channel. In *Proceedings of the 2012 ACM conference on Computer and communications security*, pages 628–637, 2012.
- [102] K. Stock. Mining location from social media: A systematic review. *Computers, Environment and Urban Systems*, 71:209–240, 2018.
- [103] L. Sweeney. k-anonymity: A model for protecting privacy. *International journal of uncertainty, fuzziness and knowledge-based systems*, 10(05):557–570, 2002.
- [104] H. Tian, M. Zhang, X. Luo, F. Liu, and Y. Qiao. Twitter user location inference based on representation learning and label propagation. In *Proceedings of The Web Conference 2020*, pages 2648–2654, 2020.
- [105] J. C. Valverde-Rebaza, M. Roche, P. Poncelet, and A. de Andrade Lopes. The role of location and social strength for friendship prediction in location-based social networks. *Information Processing & Management*, 54(4):475–489, 2018.
- [106] Y. Wang, X. Fan, X. Liu, C. Zheng, L. Chen, C. Wang, and J. Li. Unlicensed taxis detection service based on large-scale vehicles mobility data. In *2017 IEEE International Conference on Web Services (ICWS)*, pages 857–861, 2017. doi: 10.1109/ICWS.2017.106.
- [107] X. Wei, Y. Qian, C. Sun, J. Sun, and Y. Liu. A survey of location-based social networks: problems, methods, and future research directions. *GeoInformatica*, pages 1–41, 2022.
- [108] B. Wing and J. Baldrige. Simple supervised document geolocation with geodesic grids. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 955–964, 2011.
- [109] B. Wing and J. Baldrige. Hierarchical discriminative classification for text-based geolocation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 336–348, 2014.
- [110] A. G. Woodruff and C. Plaunt. Gipsy: Automated geographic indexing of text documents. *Journal of the American Society for Information Science*, 45(9):645–655, 1994.
- [111] Y. Wu, D. Lian, S. Jin, and E. Chen. Graph convolutional networks on user mobility heterogeneous graphs for social relationship inference. In *IJCAI*, pages 3898–3904, 2019.

- [112] X. Xiao, Y. Zheng, Q. Luo, and X. Xie. Finding similar users using category-based location history. In *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*, pages 442–445, 2010.
- [113] D. Xu, P. Cui, W. Zhu, and S. Yang. Find you from your friends: Graph-based residence location prediction for users in social media. In *2014 IEEE international conference on multimedia and expo (ICME)*, pages 1–6. IEEE, 2014.
- [114] Y. Yamaguchi, T. Amagasa, and H. Kitagawa. Landmark-based user location inference in social media. In *Proceedings of the first ACM conference on Online social networks*, pages 223–234, 2013.
- [115] F. Yi, H. Li, H. Wang, H. Wen, and L. Sun. Mobility intention-based relationship inference from spatiotemporal data. In *Wireless Algorithms, Systems, and Applications: 12th International Conference, WASA 2017, Guilin, China, June 19-21, 2017, Proceedings 12*, pages 871–876. Springer, 2017.
- [116] J. J.-C. Ying, E. H.-C. Lu, W.-C. Lee, T.-C. Weng, and V. S. Tseng. Mining user similarity from semantic trajectories. In *Proceedings of the 2nd acm sigspatial international workshop on location based social networks*, pages 19–26, 2010.
- [117] Q. Yuan, G. Cong, Z. Ma, A. Sun, and N. M. Thalmann. Who, where, when and what: discover spatio-temporal topics for twitter users. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 605–613, 2013.
- [118] Y. Zhang, M. Humbert, T. Rahman, C.-T. Li, J. Pang, and M. Backes. Tagvisor: A privacy advisor for sharing hashtags. In *Proceedings of the 2018 World Wide Web Conference, WWW '18*, page 287–296, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee. ISBN 9781450356398. doi: 10.1145/3178876.3186095. URL <https://doi.org/10.1145/3178876.3186095>.
- [119] W. X. Zhao, N. Zhou, W. Zhang, J.-R. Wen, S. Wang, and E. Y. Chang. A probabilistic lifestyle-based trajectory model for social strength inference from human trajectory data. *ACM Transactions on Information Systems (TOIS)*, 35(1):1–28, 2016.
- [120] Y. Zhao, M. Qiao, H. Wang, R. Zhang, D. Wang, K. Xu, and Q. Tan. Tdfi: Two-stage deep learning framework for friendship inference via multi-source information. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, pages 1981–1989. IEEE, 2019.



- [121] X. Zheng, J. Han, and A. Sun. A survey of location prediction on twitter. *IEEE Transactions on Knowledge and Data Engineering*, 30(9):1652–1671, 2018.
- [122] F. Zhou, B. Wu, Y. Yang, G. Trajcevski, K. Zhang, and T. Zhong. Vec2link: Unifying heterogeneous data for social link prediction. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 1843–1846, 2018.
- [123] W.-Y. Zhu, W.-C. Peng, L.-J. Chen, K. Zheng, and X. Zhou. Modeling user mobility for location promotion in location-based social networks. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1573–1582, 2015.
- [124] X. Zhu and Z. Ghahramani. Learning from labeled and unlabeled data with label propagation. *ProQuest Number: INFORMATION TO ALL USERS*, 2002.



## List of Figures

2.1	Generalization-Based Mechanisms. . . . .	13
2.2	Dummies-Based Mechanisms. . . . .	14
2.3	Path Confusion-Based Mechanisms. . . . .	15
2.4	Conceptual map of LPPMs. . . . .	16
3.1	Label Propagation Algorithm example. . . . .	35
3.2	Meta Decision Tree example. . . . .	41
4.1	Co-locations between friends example. . . . .	51
4.2	Trajectory example of two students that are friends. . . . .	53



## List of Tables

2.1	An example dataset with $k$ -anonymity with $k = 2$ . . . . .	11
2.2	Privacy models and LPPMs . . . . .	18
3.1	Works regarding various types of Localization through Social Media . . . .	32
3.2	Works regarding Localization through Social Network . . . . .	48

