

Improving Orienteering-based Tourist Trip Planning with Social Sensing

Fabio Persia

Faculty of Computer Science, Free University of Bozen-Bolzano, Italy

Giovanni Pilato

*Institute for High Performance Computing and Networking, National Research Council of
Italy, Palermo*

Mouzhi Ge

Faculty of Informatics, Masaryk University, Czech Republic

Paolo Bolzoni

School of Computing, Dublin City University, Ireland

Daniela D'Auria

Faculty of Computer Science, Free University of Bozen-Bolzano, Italy

Sven Helmer

Department of Informatics, University of Zurich, Switzerland

Abstract

We enhance a tourist trip planning framework based on orienteering with category constraints by adding social sensing. This allows us to customize a user's experience without putting the burden of preference elicitation on the user. We identify the interests of a user by analyzing their Tweets and then match these interests to descriptions of points of interests. For this analysis we adapt different schemes for social sensing to the needs of our orienteering context and compare them to find the most suitable approach. We show that our tech-

Email addresses: fabio.persia@unibz.it (Fabio Persia),
giovanni.pilato@icar.cnr.it (Giovanni Pilato), mouzhi.ge@muni.cz (Mouzhi Ge),
paolo.bolzoni@dcu.ie (Paolo Bolzoni), daniela.dauria@unibz.it (Daniela D'Auria),
helmer@ifi.uzh.ch (Sven Helmer)

Author Accepted Manuscript (AAM) of the paper:

Persia, F., Pilato, G., Ge, M., Bolzoni, P., D'Auria, D., & Helmer, S. (2020). Improving orienteering-based tourist trip planning with social sensing. *Future Generation Computer Systems*, 110, 931-945. <https://doi.org/10.1016/j.future.2019.10.028>

nique is fast enough for use in real-time dynamic settings and also has a higher accuracy compared to previous approaches. Additionally, we integrate a more efficient algorithm for solving the orienteering problem, boosting the overall performance and utility of our framework further, as demonstrated by the positive user satisfaction received by real users.

Keywords: social sensing, orienteering, semantic mapping, semantic similarity

1 Introduction

In tourist trip planning, it is generally not useful to determine shortest paths from a source location to a destination location. Instead, it is more important to discover routes covering the most attractive points of interests (POIs), ideally matching the personal tastes of a tourist [1]. A common and well-known way to model tourist trip planning is by mapping it onto an *orienteering problem*. The orienteering problem (OP) considers the route network as an edge-weighted graph with a weight and a score for each node. The goal is to discover a path from a starting node to a destination node that maximizes the total score while staying within a certain time budget. Having its roots in the longest path problem, OP is an NP-hard problem, which means that computing exact solutions is only possible for small problem instances.

There are more efficient algorithms for OP based on heuristics or approximations, which provide good answers in a fraction of the time it needs to compute the optimal answer. In the context of tourist trip planning this is fine, as tourists are already content with a route close to the optimal one. However, these approaches assume that the user preferences are given as input parameters to an algorithm, but do not specify where these preferences come from. In the case of category-based OP, a user would have to specify which categories of POIs they like and to which degree they like them, otherwise an algorithm cannot personalize the route.

The proliferation of numerous online social networks, such as *Facebook*,

23 *LinkedIn*, and *Twitter*, has now made it possible to incorporate social elements
24 into a variety of data-centric applications [2]. The sheer amount of information
25 found in social networks is huge and on top of that it is heterogeneous, often
26 containing complex linkage structures. *Social sensing* is an approach to analyze
27 user-generated content in social networks in such a way as to make it usable
28 for a range of applications [3]. We propose a method based on social sensing
29 to obtain personalized POI scores without explicitly asking users about their
30 specific tastes.

31 More specifically, we look at the Twitter posts of users and compute semantic
32 distances between a user profile and the description of POIs. In this way we
33 obtain a list of POIs together with personalized scores. After normalizing these
34 scores, they can be fed directly into an OP solver together with the route graph
35 to generate a custom-tailored trip without having to modify the original OP
36 algorithm. In summary, we focus on the following aspects and make a number
37 of contributions:

- 38 • We propose an accurate social sensing technique for orienteering by com-
39 bining schemes such as *Latent Dirichlet Allocation (LDA)*, *Latent Seman-*
40 *tic Analysis (LSA)*, the *spaCy*¹ tool, and the *RAKE* algorithm [4];
- 41 • For solving the associated OP, we integrate an efficient anytime algorithm
42 that was developed recently [5], making it possible to show incrementally
43 improving solutions to users;
- 44 • We evaluate our approach empirically using real-world data sets of Tweets
45 and descriptions of POIs in Italian; we opted for Italian rather than En-
46 glish due to two reasons; first, Bolzano being an Italian city, most of the
47 descriptions of its POIs are in Italian; second, as there are significantly
48 fewer applications in Italian, a framework in this language offers more
49 novelty;

¹<https://spacy.io/>

- 50 • We demonstrate the effectiveness of our approach in terms of user satis-
51 faction via a user study in which users get to compare and choose from
52 solutions generated with and without social sensing;

53 The remainder of the paper is organized as follows: in Section 2 we provide
54 a motivating example to illustrate the intuition of our work, whereas in Section
55 3 we review the technologies that are used in this paper. In Section 4 the overall
56 process is presented, and the experimental evaluation is exhibited in Section 5.
57 Eventually, Section 6 concludes the paper and outlines the future work.

58 2. Motivating Example

59 We now introduce a motivating example to illustrate the usefulness of social
60 sensing for orienteering. We take POIs from the city of Bolzano: Figure 1 shows
61 a graph with a source node s (*Waltherplatz*), a destination d (*Via Rovigo*), and
62 three POIs, p_1 (*Cathedral of Bolzano*), p_2 (*Christ of the King Church*), and
63 p_3 (*Dr. Friedrich Tessmann Library*), along with their distances expressed in
64 minutes. We simplify the graph to keep it readable: all visiting times are
65 supposed to be 60 minutes, all POIs belong to the same category, and we also
66 omit some edges.

67 Let us assume that Elon Musk (co-founder, CEO, and product architect of
68 Tesla, Inc., as well as one of the top-200 most followed twitter accounts) would
69 like to spend 3 hours (i.e., 180 minutes) in Bolzano, starting from Waltherplatz
70 (s) and ultimately arriving at a friend’s apartment (located at d). The algorithm
71 in [6], without social sensing, would suggest the itinerary $I_1 = \langle s, p_1, p_2, d \rangle$,
72 assuming all the scores to be the same. In fact, if all the other parameters are
73 the same, the algorithm suggests the itinerary with the best centrality, that is
74 the one containing the nodes with the shortest distance to s and d (for more
75 details on the algorithm, see Section 4.6).

76 On the other hand, the new algorithm including *social sensing* extracts from
77 Elon Musk’s tweets interesting topics such as “project”, “Tesla”, and “physics”,
78 that implies that p_3 (the *Dr. Friedrich Tessmann Library*, one of the most

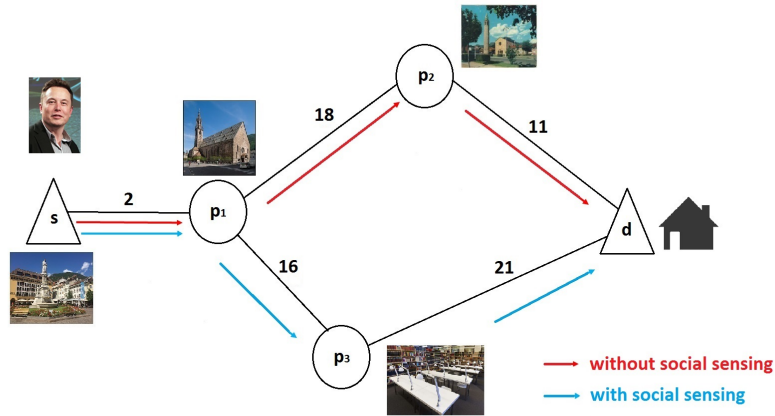


Figure 1: Route Orienteering for Elon Musk

79 modern and well-stocked libraries in Italy), is implicitly ranked higher than p_2
 80 by Elon Musk. As a result, its score will be higher than the one of p_2 , and
 81 consequently the algorithm suggests the itinerary $I_2 = \langle s, p_1, p_3, d \rangle$; even if it is
 82 longer than I_1 , it is likely to be appreciated more by this particular twitter user.
 83 So, the most interesting path for him in Bolzano would start at Waltherplatz,
 84 go to the Cathedral in two minutes, spend an hour there, then arrive at the Dr.
 85 Friedrich Tessmann Library sixteen minutes later, stay an hour there, and reach
 86 the friend's apartment in Via Rovigo 21 minutes later. This path has a total
 87 length of 159 minutes, leaving more than 20 minutes to the time constraint of
 88 three hours.

89 3. Related Work

90 The *orienteering problem* was initially defined by Tsiligrides in [7]. Since it is
 91 an NP-hard problem, algorithms computing the exact solutions via branch and
 92 bound or dynamic programming techniques [8, 9] can only solve small instances

93 of the problem. As a consequence, most of the used approximation algorithms
94 and heuristics adopt a two-step approach of partial path construction [10, 7]
95 and (partial) path improvement [11, 12, 13]. Additionally, other approaches in
96 the literature exploit meta-heuristics such as genetic algorithms [14] and neural
97 networks [15]. There are also a few approaches that consider the orienteering
98 problem generalized by categories. These approaches are based on clustering the
99 nodes [16], applying an optimized best-first search [5], and using probabilistic
100 methods [6]. However, none of these techniques consider using social sensing for
101 eliciting user preferences.

102 In order to integrate social network information into the solution of an orien-
103 teering problem, the first step is to process the *Online Social Network (OSNs)*
104 data. There are two general categories of methods that can extract, collect,
105 and mine the information from *OSN* data: *machine learning* approaches and
106 *non-machine learning* approaches. An example of a *machine learning* approach
107 is [17], which uses a range of features including language, hashtags, and geo-
108 graphical locations of their social ties to infer a user’s nationality on Twitter;
109 thus, it can be classified as a method based on *user profile*. Furthermore, Bent-
110 wood [18] uses a topology of influence to quantify the impact of social media,
111 while Tinati et al. [19] uses a social network graph based on retweet behavior
112 to identify key players in a conversation. They describe a new approach to clas-
113 sify users within the Twitter service and use a topology of influence. Moreover,
114 D’Avanzo et al. [20] introduce a pipeline, implemented as a web service, to al-
115 low decision makers to monitor Twitter’s sentiment regarding trends based on
116 Google Trends. It also enables users to choose geographic areas for their mon-
117 itors. In addition to the positive/negative sentiments on Google Trends, the
118 pipeline offers the ability to view, on the same dashboard, the emotions that
119 Google Trends triggers in the Twitter population. For the *non-machine learn-*
120 *ing* approaches, Gianvecchio et al. [21] propose an *entropy-based* classification
121 method to identify the chat bots in Yahoo Chat, while Sarwar et al. [22] exploit
122 *user-based* collaborative filtering to predict a test user’s interest in a test item
123 based on rating information from similar user profiles.

124 We now look at some well-known techniques and tools used by us for sensing
 125 the interests of user from their social network profiles. More specifically, in
 126 subsection 3.1 we show the approach built on the LSA-based Semantic Space,
 127 in subsection 3.2 the one built on the LDA-based Topics Space, in subsection 3.3
 128 the one based on an effective NLP tool named spaCy, and in subsection 3.4 the
 129 one based on the RAKE algorithm.

130 3.1. LSA-based Semantic Space

131 The Latent Semantic Analysis (LSA) technique is a well-known method use-
 132 ful for obtaining a vector encoding of words and their semantics [23]. LSA has
 133 been widely used for simulating different human cognitive phenomena [24][25].
 134 Given a set of M terms and N documents, the LSA technique is based on a
 135 term-document co-occurrence $M \times N$ matrix \mathbf{A} , whose generic element $[a_{i,j}]$ is a
 136 function of the number of times a term occurs in a document. In particular, we
 137 have considered the classical “Term frequency - Inverse document frequency”
 138 *Tf-Idf* approach for encoding a generic element $[a_{i,j}]$ of \mathbf{A} . Given a word w_i
 139 and a document d_j , let $f_{i,j}$ be the number of occurrences of w_i in d_j and let n_{w_i}
 140 be the number of times the word w_i appears at least once in a document of the
 141 collection. Then, a generic element $[a_{i,j}]$ of the matrix \mathbf{A} is:

$$[a_{i,j}] = (1 + \log(f_{i,j})) \cdot \log\left(\frac{N}{n_{w_i}}\right) \quad (1)$$

142 The matrix \mathbf{A} can be decomposed into the product of three matrices, \mathbf{U} , $\mathbf{\Sigma}$,
 143 and \mathbf{V} by exploiting a factorization procedure known in literature as “Singular
 144 Value Decomposition” (SVD)[23].

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (2)$$

145 \mathbf{U} is a column-orthonormal $M \times N$ matrix, \mathbf{V} is a column-orthonormal $N \times N$
 146 matrix and $\mathbf{\Sigma}$ is a $N \times N$ diagonal matrix. The elements of $\mathbf{\Sigma}$ are named
 147 “singular values” of \mathbf{A} . The columns of \mathbf{U} are named left singular vectors of
 148 \mathbf{A} and the columns of \mathbf{V} are named right singular vectors of \mathbf{A} . It may be
 149 assumed, without loss of generality, that the singular values of \mathbf{A} are ranked in

150 decreasing order, since Singular Value Decomposition also holds if we perform
 151 permutations of singular values along with the corresponding rows in \mathbf{U} and \mathbf{V}
 152 [25].

153 Let R be an integer with $0 < R < N$, and let \mathbf{U}_R be the $M \times R$ matrix extracted
 154 from \mathbf{U} by neglecting the last $N - R$ columns, let Σ_R be the matrix extracted
 155 from Σ by neglecting the last $N - R$ rows and the last $N - R$ columns; let \mathbf{V}_R
 156 be the $N \times R$ matrix extracted from \mathbf{V} by neglecting the last $N - R$ columns.
 157 Then, it can be shown that:

$$\mathbf{A}_R = \mathbf{U}_R \Sigma_R \mathbf{V}_R^T \quad (3)$$

158 \mathbf{A}_R is an $M \times N$ matrix of rank R , and it is the best rank R approximation of
 159 \mathbf{A} , among the $M \times N$ matrices, with respect to the Euclidean metric. The i -th
 160 row of the matrix \mathbf{U}_R is the numerical encoding of the i -th word. The columns
 161 of the \mathbf{U}_R matrix represent the R independent dimensions of the induced se-
 162 mantic space. Moreover, it is possible to extract the words which best describe
 163 these independent dimensions by computing a score $Sc(w_i)$ for each word w_i by
 164 applying the formula [25, 26]:

$$Sc(w_i) = \sum_{j=1}^R \sigma_{ij} |u_{ij}| \quad (4)$$

165 with $\sigma_{ij} \in \Sigma_R$ and $u_{ij} \in U_R$.

166 3.2. LDA-based Topics Space

167 Latent Dirichlet Allocation (LDA) is an unsupervised method for computing
 168 a Bayesian probabilistic model of text corpora, with the aim of finding topics
 169 in documents [27].² In a nutshell, it assumes that the documents in a corpora
 170 were created using a generative process (for a gentle introduction to LDA and
 171 generative probabilistic topic models, see [29]). Put simply, we assume that
 172 every document follows a certain topic distribution (e.g. 70% topic A and

²LDA extends the probabilistic Latent Semantic Analysis (pLSA) approach [28], which, even though its name looks similar to LSA, follows a completely different procedure.

173 30% topic B). Furthermore, a topic is specified as a distribution over a fixed
 174 vocabulary (i.e., the words contained in a document collection). Generating a
 175 document is split into two stages: first we (randomly) choose a distribution over
 176 topics for the document and then every word is randomly determined by the
 177 distribution over the vocabulary while following the topic distribution of the
 178 document. The main issue with this process is that we know neither the topic
 179 distributions of the documents nor the vocabulary distributions of the topics,
 180 these are hidden variables. LDA turns this around by searching for the most
 181 probable hidden parameters that could have generated an observed document
 182 collection. More formally, this is done in the following way. The number of
 183 topics, K , is fixed beforehand and every topic $\beta_k \in \beta_{1:K}$ is a distribution over
 184 the vocabulary. Now we turn to the document collection, which contains a total
 185 of D documents. The percentage of topic k in document d is expressed by $\theta_{d,k}$,
 186 θ_d defines the percentages of all the topics in this document. The distribution
 187 $z_{d,n}$ determines which topic is assigned to the word n in document d , z_d stands
 188 for all these assignments in d . Last, but not least, $w_{d,n}$ represents the n -th
 189 word in document d , w_d denotes all the observed words in d . Putting all this
 190 together, we can now describe the whole process of generating documents using
 191 a joint probability distribution:

$$\begin{aligned}
 & p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D}) = \\
 & \prod_{i=1}^K p(\beta_i) \prod_{d=1}^D p(\theta_d) \left(\prod_{n=1}^N p(z_{d,n} | \theta_d) p(w_{d,n} | \beta_{1:K}, z_{d,n}) \right) \quad (5)
 \end{aligned}$$

192 On a technical note: the different distributions above are modeled using sym-
 193 metric Dirichlet distributions. The parameters β , θ , and z are called latent
 194 variables, as they are not directly observable. Given the observations (in the
 195 form of word frequencies in the documents), we reverse-engineer the generative
 196 process by computing the posterior distributions of the latent variables, which
 197 boils down to solving

$$\mathbf{P}(\beta_{1:K}, \theta_{1:D}, \mathbf{z}_{1:D} | \mathbf{w}_{1:D}) = \frac{\mathbf{P}(\beta_{1:K}, \theta_{1:D}, \mathbf{z}_{1:D}, \mathbf{w}_{1:D})}{\mathbf{P}(\mathbf{w}_{1:D})} \quad (6)$$

198 The process of solving this equation is intractable from a computational point of
199 view. However, its solution can be approximated through variational inference
200 and Gibbs sampling [30, 27, 31].

201 3.3. *SpaCy Tool*

202 For this method, we used two natural language processing (NLP) tools,
203 namely fastText [32] and spaCy [33], expanding and improving the spaCy base-
204 line approach we applied in [34]. Our technique is composed of two steps:
205 preprocessing documents with the help of fastText and comparing them via
206 spaCy.

207 Before starting with the preprocessing, we remove all links to webpages, ref-
208 erences to other users, and hashtags from our collection of Tweets. We then feed
209 the words contained in the Tweets and the POI descriptions into fastText, which
210 computes distributed word representations, or word vectors, for them [35, 36].
211 This is a popular technique in machine learning for uncovering subsymbolic
212 meanings, such as word analogies. We utilized a pre-trained word vector encod-
213 ing for Italian provided by fastText [32], which was trained on Common Crawl
214 and Wikipedia. As reported by fastText, the parameters used for this training
215 are as follows: 300-dimensional vectors, a continuous bag of words (CBOW)
216 model with position weights, character n-grams of length five, a window size of
217 five, and ten negative examples [32, 36].

218 We use a simple strategy for determining the similarity between Twitter
219 profiles and POI descriptions. Assuming that each POI is described by (at least)
220 one sentence, we compute the similarity between the word vectors describing
221 the Tweets and the word vector of the sentence describing a POI via the cosine
222 measure. If the similarity between a Tweet and a POI is greater than a threshold
223 θ , the Tweet matches the POI (the value of θ is set to 0.85, which was determined
224 experimentally). For measuring the similarity of the whole Twitter account to a
225 POI, we take the number of matching Tweets and divide it by the total number
226 of Tweets for this account. If a POI is described by more than one sentence, we
227 only consider the sentence with the highest similarity. We repeat this procedure

228 for every POI, obtaining an ordered list of POIs according to the semantic
229 similarity with the user posts.

230 In our previous work we showed that a social sensing approach based on LDA
231 and LSA is able to achieve an accuracy similar to a simpler baseline version of
232 the spaCy approach, but at a much faster rate [34]. Here, we build on top of this
233 by optimizing the spaCy approach (see scheme 5 in Section 4.3.3), decreasing
234 the number of operations per user it needs.

235 3.4. Rapid Automatic Keyword Extraction (RAKE)

236 The RAKE algorithm [4] has been introduced by Rose et al. with the in-
237 tention to get a general and efficient algorithm for extracting keywords from
238 documents. Keywords are defined as sequences of one or more words capable of
239 supplying a concise representation of the subject of a document [4].

240 RAKE is based on the observation that keywords often contains multiple
241 words, but they seldom contain stop words in the middle. In addition to the
242 documents, the method takes as input a set of of stop words, e.g. for English
243 these are words such as “the”, “of”, or “and”, and a set of phrase and words
244 delimiters like blanks or punctuation marks. The stop words and delimiters
245 are used to partition the text into word sequences, with each sequence being
246 considered a candidate keyword.

247 The algorithm is unsupervised, domain-, and language-independent. In fact,
248 the original paper [4] also describes a method to automatically generate a list
249 of stop words, starting with a set of documents.

250 Once the list of the candidate keywords has been identified, each word w_i
251 present in the text is scored by using the function:

$$s_{w_i} = \frac{degree(w_i)}{occ(w_i)} \quad (7)$$

252 where $occ(w_i)$ is the number of occurrences of w_i in the text, and $degree(w_i)$ is
253 the number of times that the word w_i co-occurs with other words in a candidate
254 keyword sequence. A higher value of word $degree(w_i)$ means that w_i occurs in

255 long candidate keywords. The index i ranges from one to the number of unique
256 words occurring in the text.

257 Once the score of each word in the text has been computed, a grade is
258 associated with each candidate keyword by summing the scores of the words
259 that constitute it. The candidate keywords are then sorted by descending value
260 of their grades and the top third are selected as keywords of the text [4].

261 The algorithm can also be enabled to detect keywords with stop words in the
262 middle. In this case, once the document in question is parsed and a set of words
263 is found, the algorithm looks for pairs of possible keywords with a stop word
264 in the middle that appear at least twice in the same document and includes
265 them among the list of keyword candidates. This extra step allows us to obtain
266 keywords like “*set of natural numbers*”, where the presence of the stop word
267 “*of*” would have prevented us from finding it.

268 Experiments reported in [4] show that RAKE is an efficient algorithm, while
269 at the same time achieving a good ratio of precision and recall. In this paper,
270 RAKE has been used as one of the schemes to extract keywords from user
271 tweets.

272 4. Orienteering based on Social Sensing

273 In this section, we show the overall process we followed for building an ori-
274 enteering framework supported by social-sensing (Figure 2). It consists of six
275 layers: *OSN crawler*, *tweets extraction and preprocessing*, *semantic space map-*
276 *ping*, *semantic similarity matrix building*, *ranking computation*, and *integration*
277 *with the orienteering approach*.

278 The output of the OSN crawler consists of a set Λ of collected data related
279 to user sessions on a particular OSN. As a result, in this context we collect all
280 the *user tweets*, which are further processed in the following steps (Figure 2) of
281 the pipeline, in order to detect the *user interests* in *POIs*.

282 Then, the *tweets extraction and preprocessing step* aims at cleaning the text
283 of the tweets before further processing. Afterwards, the *semantic space mapping*

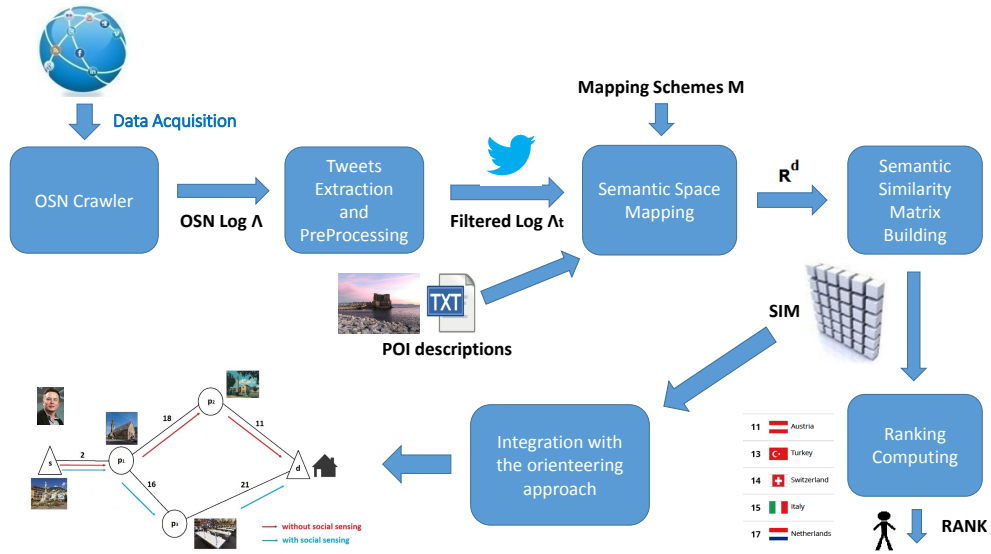


Figure 2: The Overall Process

284 phase is applied to map both the filtered tweets and the POI descriptions into
 285 a joint semantic space \mathbf{R}^d ; the details of this method are explained later on.

286 With the joint semantic space, the overall *semantic similarity matrix* SIM
 287 between *users* and *POIs* is computed; a generic element $s_{ij} \in SIM$ stands for
 288 the *user interest* of user i in POI j , computed by analyzing the similarity of
 289 the *user tweets* and the *POI descriptions*. Afterwards, two further processing
 290 steps are possible; on the one hand, a specific *RANK* matrix holding the user
 291 interests in the POIs ordered by descending values is extracted; this is done
 292 to ease the process of extracting the *Top-k* list of preferences for each user,
 293 where $0 < k < n$, n being the total number of considered POIs. The *RANK*
 294 matrix can be used in an *offline* setting as an input for applications such as POI
 295 visualization or recommender systems.

296 On the other hand, the SIM matrix can be exploited for customizing the
 297 scores of the nodes as used in an algorithm for solving the orienteering problem.

298 We provide customized paths that also take into account the *user interest* in
299 the POIs, implicitly extracted from a user’s tweets. As a result, the score of each
300 POI (Figure 2) can be considered as a function of the similarities among the
301 topics of the user tweets and the topics of the POI descriptions. The detailed
302 procedure is described in Subsection 4.6.

303 4.1. OSN Crawler

304 We start off by considering a finite set of atomic user actions in an OSN along
305 the lines of [37] (Table 1). We assume that the set of observable actions in the
306 domain is finite, since we deal with sequences of time-stamped events. Each
307 action can be associated with a different action symbol in S , the user to whom
308 it belongs, a related timestamp, and other useful attributes for characterizing
309 actions. In this paper we analyze *Twitter posts* and focus on the “e” symbol
310 of Table 1. Thus, we store the *content of tweets* as additional information for
311 characterizing the action. We base the definition of our *OSN Log* on the one
312 provided in [37, 38, 39].

313 *OSN Log*: A *log entry* is a tuple $\lambda = (s, user, ts, att_1, \dots, att_k)$, where $\lambda.s \in S$
314 is the action symbol associated with the event, $\lambda.user$ is the user with which
315 the action is associated, $\lambda.ts$ is the related timestamp and $\lambda.att_1 \dots \lambda.att_k$ are
316 other attributes optionally associated with the observation (e.g. the content of
317 a tweet, or an IP address). An *OSN Log* Λ is then a finite sequence of log entries
318 λ_i ordered by timestamp.

319 In this context, the task of the *OSN Crawler* is the collection of real datasets,
320 such as the one described in subsection 5.2.1, using the Twitter APIs with the
321 default access level.

322 4.2. Tweets extraction and preprocessing

323 Given a user Twitter ID, we extract their tweets. For our approach, we
324 only use the tweet text content, which is preprocessed before being utilized to
325 build a data-driven conceptual space. Stop-words are filtered out, and links
326 are removed before processing the text, since they often hide off-topic posts or

Table 1: Action/Symbol

Symbol	Action
a	login
b	message received
c	message sent
d	friend approved
e	status wall post
f	link shared story
g	video shared story
h	status update
i	added picture
j	published link
k	pictured shared story
l	youtube video shared
m	youtube created story
n	link app created story
o	tagged in a picture
p	mobile status update
q	checkin
r	likes a page
s	logout

327 even spam. Abnormal sequences of characters are discarded and the twitter
328 text is then lemmatized. In this context the *OSN Crawler* can be used either
329 for retrieving tweets satisfying a query composed of keywords, or to download
330 the most recent tweets of a given twitter user ID. The first step of the followed
331 procedure is the text extraction. More specifically, for each Twitter account x
332 we retrieve all the most recent user tweets via the Twitter APIs and then store
333 them. As a result, in our case study the maximum number of tweets detected
334 per user is about 3200. We also detect the user’s followers and friends, and save
335 them for possible future use. Afterwards, we identify the language of the tweets;
336 if less than 60% of the user tweets are in Italian or the user tweets in Italian
337 are less than 20, we ignore the Twitter account. Otherwise, we extract from the
338 tweets only the textual content (getting rid of all the other information).

339 4.3. Semantic scoring of POIs according to the user interests

340 Given the user Twitter ID, we have set up a set of approaches that can help
341 in finding a semantic similarity between the user tweets and the natural lan-
342 guage description of each POI. In particular, we defined six different schemes
343 exploiting a semantic space induced from user tweets by Latent Semantic Analy-
344 sis (LSA) [23][40], a topic space induced by the user tweets from Latent Dirichlet
345 Allocation (LDA) [27], the *spaCy* tool, and the RAKE algorithm [4]. As dis-
346 cussed in Section 3, LSA, LDA, *spaCy*, and RAKE algorithm are techniques
347 that can also be used for profiling users and to extract the main topics in tweets
348 [26][41]. We explain all the different approaches for getting the semantic scoring
349 of POIs according to the user interests in Sections 4.3.1 to 4.3.4.

350 4.3.1. Semantic scoring schemes 1 and 2

351 The weights identified by the schemes with numbers 1 and 2 are respectively:

- 352 1. the *maximum* value of Jaccard similarity between topics induced by using
353 LSA and LDA from the tweets of the user, and the topic descriptions;
- 354 2. the *average* value of Jaccard similarity between topics induced by using
355 LSA and LDA from the tweets of the user, and the topic descriptions.

356 In order to obtain these weights, we apply the following procedure. First,
 357 we remove from the textual content of the tweets any links, references to other
 358 accounts, and hashtags (since the POI descriptions do not contain any of such
 359 elements); then, we also remove stopwords. Consequently, on such filtered tweets
 360 of the i -th user (that we identify as set Λ_{ti}) we execute LSA, thus obtaining the
 361 first 20 *conceptual axes* along with their descriptive words (10 per axis). Then,
 362 we also apply LDA on Λ_{ti} , thus obtaining 20 topics along with their descriptive
 363 words (10 per topic). As a result, by applying *LSA* and *LDA* on Λ_{ti} , we obtain
 364 a file containing 40 rows; each of them has 10 words, describing the possible
 365 basic user interests (U_I). On the other hand, we remove the stopwords from
 366 the POI descriptions; for each descriptive sentence of a POI, we compute the
 367 Jaccard similarity (Intersection/Union) with each of the 40 rows of the user as
 368 well as the maximum value obtained for each descriptive sentence, keeping also
 369 track of all the Jaccard values. More formally, we define:

- 370 • n_s as the number of sentences of the POI description;
- 371 • n_{U_i} as the total number of *conceptual* descriptions of user interests (40 in
 372 this case).

373 Consequently, the computation of the Jaccard similarities implies the con-
 374 struction of the J matrix; more specifically, in Equation 8 the *Jaccard* function
 375 is defined on the domain $\{1, \dots, n_s\} \times \{1, \dots, n_{U_i}\}$, and has the J matrix as co-
 376 domain.

$$Jaccard : \{1, \dots, n_s\} \times \{1, \dots, n_{U_i}\} \rightarrow J \quad (8)$$

377 Eventually, for each tuple ($user, POI$), we compute:

- 378 • the average of the maximum values (*weight 1*, Equation 9);
- 379 • the average of all the Jaccard values (*weight 1*, Equation 10).

$$w_{ij}^{(1)} = avg\{\max_k\{J_{k*j}\}\}, \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n; \forall k = 1, \dots, n_s; \quad (9)$$

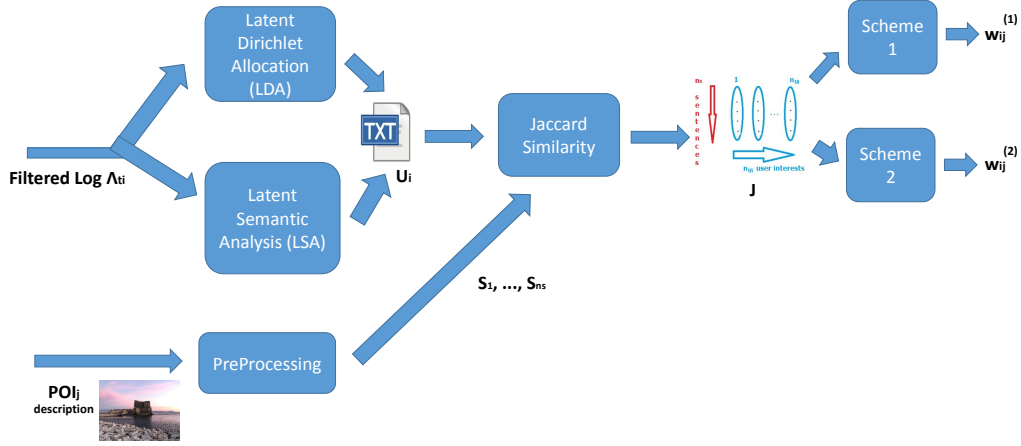


Figure 3: Scheme 1 and 2, for each tuple $(user_i, POI_j)$

$$w_{ij}^{(2)} = avg\{J_{kl}\}, \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n; \forall k = 1, \dots, n_s; \forall l = 1, \dots, n_{U_i} \quad (10)$$

380 4.3.2. Semantic scoring schemes 3 and 4

381 The weights identified by the schemes with numbers 3 and 4 are, respectively:

- 382 1. the *maximum* value of Jaccard similarity between the user tweets and the
 383 topic descriptions;
- 384 2. the *average* value of Jaccard similarity between the user tweets and the
 385 topic descriptions.

386 The procedure to obtain these weights is described in the following. First, we
 387 remove from the textual content of the tweets links, references to other accounts,
 388 hashtags (since in the POI descriptions there are no such elements), and the
 389 stopwords. Then, we also remove the stopwords from the POI descriptions.

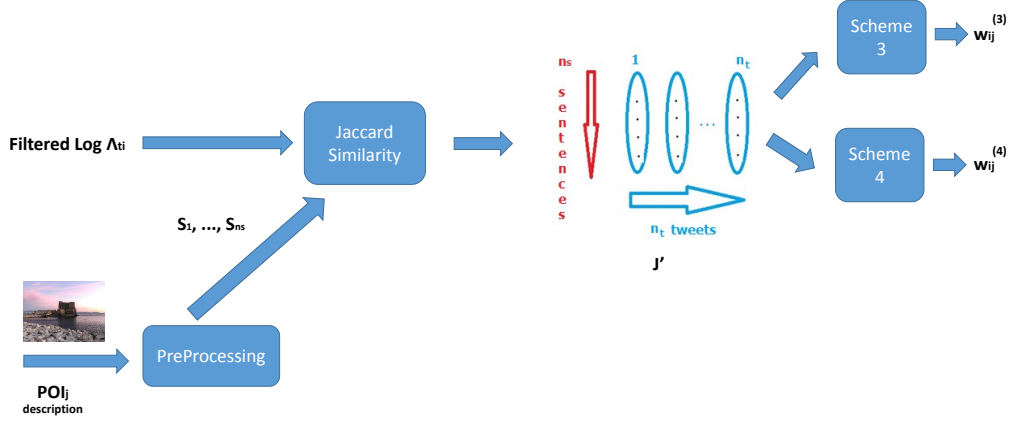


Figure 4: Scheme 3 and 4, for each tuple $(user_i, POI_j)$

390 Consequently, for each descriptive sentence of the POI we compute the Jaccard
 391 similarity (Intersection/Union) with every user tweet. Formally, assuming n_t to
 392 be the number of user tweets, we build the following matrix J' .

$$Jaccard : \{1, \dots, n_s\} \times \{1, \dots, n_t\} \rightarrow J' \quad (11)$$

393 Eventually, we compute $w^{(3)}$ (Equation 12) and $w^{(4)}$ (Equation 13) for each
 394 tuple $(user, POI)$ as follows:

$$w_{ij}^{(3)} = avg\{\max_k\{J'_{k*}\}\}, \forall i = 1, \dots, m; \forall j = 1, \dots, n; \forall k = 1, \dots, n_s \quad (12)$$

$$w_{ij}^{(4)} = avg\{J'_{kl}\}, \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n; \forall k = 1, \dots, n_s; \forall l = 1, \dots, n_t \quad (13)$$

395 4.3.3. *Semantic scoring scheme 5*

396 The weight identified by scheme 5 is explained in the following:

- 397 1. all the tweets of a user taken together are treated as a single document,
398 whereas every description of a POI forms its own document. All docu-
399 ments are encoded with fastText and we compute their similarity with the
400 help of the cosine measure provided by spaCy.

401 More specifically, we follow the listed procedure:

- 402 1. we encode the concatenated tweets of the user i as a vector with fastText,
403 exploiting the “it_vector_wiki_lg” model, obtained from the Italian version
404 of Wikipedia. As a result, we get a single vector (enc_tweets) in a 300-
405 dimensional space, which is the encoding of all the tweets of user i ;
- 406 2. we follow the same procedure as the one explained in step 1 for encoding
407 each POI description. Consequently, we obtain the matrix enc_POIs ,
408 whose generic column j represents the encoding in a 300-dimensional space
409 of the j -th POI;
- 410 3. we compute for each Twitter user the semantic distance (the *cosine*) be-
411 tween all the user tweets (as a single document) and the POI descriptions.
412 More formally, remembering that m is the total number of users, and n
413 the total number of POIs, we define the weight 5 ($w^{(5)}$) as follows:

$$w_{ij}^{(5)} = \cos(enc_tweets_{i*}, enc_POIs_{*j}) \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n \quad (14)$$

414 We can sort these weights by similarity in descending order, obtaining the
415 final ranking.

416 4.3.4. *Semantic scoring scheme 6*

417 The weight identified by scheme 6 is defined by

- 418 • the similarity between descriptive keywords extracted by the RAKE algo-
419 rithm and the description of the POI.

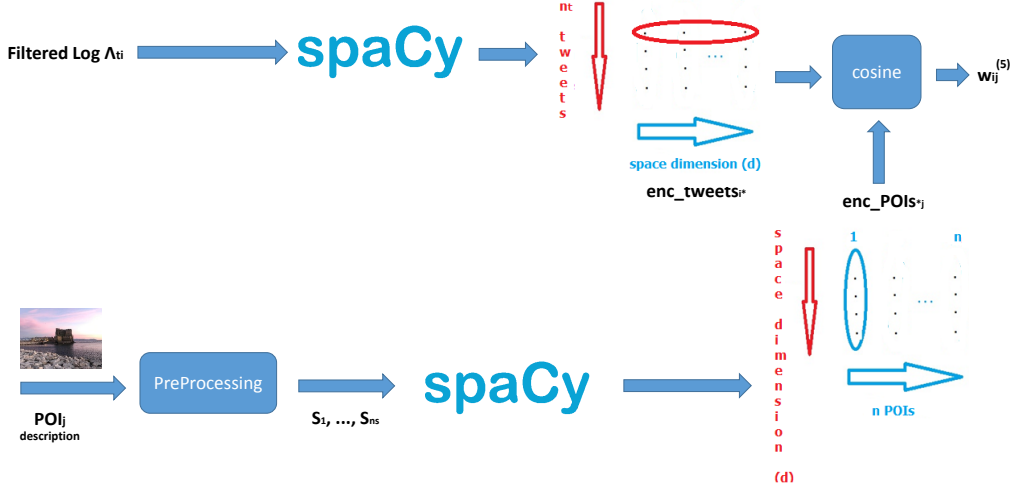


Figure 5: Scheme 5, for each tuple $(user_i, POI_j)$

420 We proceed as follows.

- 421 1. We compute the enc_POIs matrix as described in subsection 4.3.3;
- 422 2. we use the RAKE algorithm [4] to extract the n_{key} keyphrases (each of
- 423 them identified by key_k in Equation 15) from the whole set of tweets of
- 424 the i -th user;
- 425 3. each of the keyphrases is associated with a $score$ ($score_k$);
- 426 4. eventually, we build weight 6 for each tuple $(user_i, POI_j)$ as follows:

$$w_{ij}^{(6)} = \sum_{k=1}^{n_{key}} score_k \cdot \cos(key_k, enc_POIs_{*j}) \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n \quad (15)$$

427 4.3.5. Final remarks about semantic scoring schemes

428 As shown in the previous sections, all the described semantic scoring schemes
 429 map the filtered user tweets as well as the POI descriptions into \mathbf{R}^2 , i.e., $d = 2$

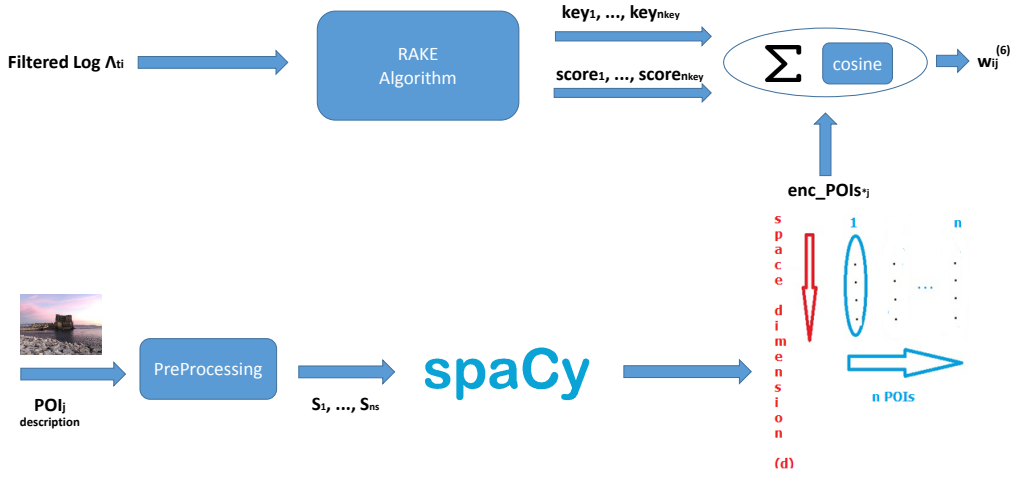


Figure 6: Scheme 6, for each tuple $(user_i, POI_j)$

430 (Figure 2). However, possible future semantic scoring schemes mapping into
 431 \mathbf{R}^d , $d > 2$, could also be integrated into the framework.

432 4.4. Semantic Similarity Matrix Building

433 In this step we combine the weights of a scoring scheme, which we discussed
 434 earlier, into a similarity *SIM*; more specifically, a generic item of the *SIM*
 435 matrix stands for the interest of user i in POI j , whereby the higher the value,
 436 the greater the interest. As a result, we define a similarity matrix for each of
 437 the defined schemes as follows:

$$SIM_{ij}^{(k)} = w_{ij}^{(k)} \quad \forall i = 1, \dots, m; \forall j = 1, \dots, n; k = 1, \dots, 6 \quad (16)$$

438 In Subsection 4.5 we show how to use the ranking computing functionality
 439 in an offline setting, whereas in Subsection 4.6 we demonstrate how to use the
 440 similarity matrix to compute customized node scores.

441 *4.5. Ranking Computing*

442 Now we show how to build a *RANK* matrix starting from a *SIM* matrix.
443 The *RANK* matrix is useful as input for applications such as POI visualiza-
444 tion, recommender systems, or data analytics; in fact, in these contexts it is
445 convenient to use a matrix whose rows are sorted in a descending order, thus
446 explicitly listing the most popular POIs in the first positions. More formally,
447 we first define the *sort* function (Equation 17), which orders a vector v in de-
448 scending order.

$$\text{sort} : \{v_1, \dots, v_m\} \rightarrow \{v'_1, \dots, v'_m\}, \mid v'_i \geq v'_{i+1}, \forall i = 1, \dots, m - 1; \quad (17)$$

449 Consequently, Equation 18 defines the *RANK* matrix as

$$RANK_{i*} = \text{sort}\{SIM_{i*}\} \quad \forall i = 1, \dots, m; \quad (18)$$

450 whose generic items stand for the interest of the user i in the POI j , sorted in
451 descending order. Clearly, for each generic item $RANK_{ij}$ we also keep track
452 of the *ID* of the related POI_j . However, since the complexity of building the
453 *RANK* matrix is $O(m \cdot n \cdot \log(n))$, it is better to exploit this functionality in an
454 offline setting (where the matrix has not to be computed too frequently), rather
455 than in a real-time setting.

456 *4.6. Integration with the Orienteering Approach*

457 Before describing the algorithm for solving the orienteering problem (OP)
458 used in [5], and the way we combine it with social sensing, we define the formal
459 model that we utilized for tourist trip planning. Since there is no obligation
460 to visit all the nodes in the OP, it is also called *Selective Traveling Salesman*
461 *Problem (STSP)*.

462 We assume a set P of points of interest (POIs) p_i , $1 \leq i \leq n$. The POIs,
463 together with a starting and a destination node, denoted by s and d , respectively,
464 are connected by a complete, metric, weighted, undirected graph $G = (P \cup$

465 $\{s, d\}, E)$, whose edges, $e_l \in E = \{(x, y) | x, y \in P \cup \{s, d\}\}$, connect them. Each
466 edge e_l has a cost $c(p_i, p_j)$ that signifies the duration of the trip from p_i to p_j ,
467 while every node $p_i \in P$ has a cost $c(p_i)$ that denotes its visiting time. Each POI
468 belongs to a certain category, such as *museums*, *restaurants*, or *galleries*. The
469 set of n_{cat} categories is denoted by K and each POI p_i belongs to exactly one
470 category $\kappa_k \in K, 1 \leq k \leq n_{cat}$. Given a POI p_i , $cat(p_i)$ denotes the category
471 p_i belongs to and $score(p_i)$ denotes its score or reward, with higher values
472 indicating higher interest to the user. Finally, users have a certain maximum
473 time in their budget to complete the itinerary, denoted by t_{max} .

474 Furthermore, let P^* be the ordered set of vertices in the optimal path, and
475 define a binary variable x_{ij} taking the value 1 if and only if p_i and p_j are
476 two consecutive vertices of P^* . Additionally, let y_{ik} be another binary variable
477 taking the value 1 if and only if POI p_i belongs to category κ_k , and $max_{\kappa_k}, \forall \kappa_k \in$
478 K represents the maximum number of POIs belonging to category κ_k allowed
479 in P^* .

480 Consequently, the orienteering problem for tourist trip planning is formally
481 modeled as follows, with $1 \leq i, j \leq n$ and $1 \leq k \leq n_{cat}$:

$$\begin{aligned} & \mathbf{max} \sum_i score(p_i) \sum_j x_{ij} \\ & \mathbf{subject\ to} \end{aligned}$$

$$\sum_i x_{si} = \sum_i x_{id} = 1 \quad (19)$$

$$\sum_i x_{il} = \sum_i x_{li} \leq 1, \quad \forall l : 1 \leq l \leq n \quad (20)$$

$$\sum_i \sum_j (c(p_i) + c(p_i, p_j)) \cdot x_{ij} \leq t_{max} \quad (21)$$

$$\sum_k y_{ik} = 1, \quad \forall i : 1 \leq i \leq n \quad (22)$$

$$\sum_i \sum_j x_{ij} \cdot y_{ik} \leq max_{\kappa_k}, \quad \forall k : 1 \leq k \leq n_{cat} \quad (23)$$

$$x_{ij} \in (0, 1); \quad i \neq j$$

482 More specifically, the constraint in Equation 19 specifies that both s and d
 483 belong to P^* as the first and last POI, respectively; additionally, constraint (20)
 484 defines the flow conservation conditions for the remaining vertices, while con-
 485 straint (21) specifies the upper bound on the path length. Moreover, constraint
 486 (22) guarantees that each POI belongs to only one category, whereas constraint
 487 (23) determines the maximum number of POIs for each category allowed in a
 488 path.

489 As a result, the main goal of integrating OP with social sensing is to utilize
 490 the *similarity matrix* previously obtained (subsection 4.4) to improve the quality
 491 of itineraries. Differently from [34], in this paper we exploit the algorithm for
 492 solving the orienteering problem developed in [5] rather than the one in [6], as
 493 it is an approximation algorithm with guaranteed bounds for the quality of the
 494 solution and as an anytime algorithm it provides a solution right from the start,
 495 which is then gradually improved.

496 Since POIs can be visited only once, Bolzoni et al. [5] use sets as route
 497 representatives. Thus, each search space node comprises a set of POIs, the
 498 related scores, and a potential extra score. Each element of the search space,
 499 then, represents a solution able to visit all the POIs in its set collecting their
 500 scores. The potential extra score measures how promising a search space node
 501 is. In particular, the algorithm computes the shortest path visiting all the POIs
 502 of a search space node, with (p_l, d) being the final arc leading to the destination.
 503 We consider the remaining budget time the remaining time at p_l . The potential
 504 extra score is an upper bound of the score we can get by visiting additional
 505 unvisited POIs between p_l and d . A trivial upper bound is to add up all the
 506 POIs p so that $c(p_l, p) + c(p, d) - c(p_l, d)$ still allows us to arrive at d
 507 in time. However, with categories it is possible to have a stricter bound, since
 508 we consider only the top scoring POIs of each category satisfying the category
 509 constraints.

510 The algorithm then uses a double-ended priority queue (aka priority deque)
 511 that keeps the solution representatives in descending order of maximum po-
 512 tential score. In more detail, routes in the deque are represented as tuples
 513 comprising $\langle c, score, extra, \check{P} \rangle$ where \check{P} is the set of POIs of the representative,
 514 c is the cost of passing through and visiting all the POIs starting from s and
 515 arriving at d , $score$ is the sum of the scores of the POIs in \check{P} , and $extra$ is an
 516 upper bound on additional scores that are still achievable. The *pot* (potential)
 517 of a representative is $score + extra$.

518 For instance, let us assume the initial representative $\check{P} = \emptyset$ for an empty
 519 route from s to d , then *cost* is the distance between s and d , *score* is zero,
 520 and *extra* contains the sum of the top scoring reachable POIs of each category.
 521 Assuming the category constraints are 2 *museums* and 1 *restaurant*, then *extra*
 522 contains the scores of the top restaurant and the two top museums still in reach.
 523 The more POIs we have in \check{P} , the higher the *cost*, and the tighter the category
 524 constraint will be, allowing *extra* to be more precise. If *extra* is zero, it means
 525 it is not possible to visit any more POIs without violating a constraint. Usually,
 526 this upper bound is not a supremum because visiting some of the top scoring
 527 POIs will make it impossible to still reach others.

528 The priority deque keeps the representatives sorted by their potential score
 529 *pot* and in each step the algorithm removes the one with maximum *pot*, repre-
 530 sented by $\langle c, score, extra, \check{P} \rangle$. Based on this removed route, we generate new
 531 representatives $\langle c_i, score_i, extra_i, \check{P} \cup p_i \rangle$ where p_i is a new POI (not in \check{P}), up-
 532 dating all its components appropriately. The new representatives are put back
 533 into the priority deque only if the sum of *score* and *extra* is higher than the
 534 score of the currently best known route.

535 Each time a representative is computed, the relative route is compared with
 536 the best known route. If it is better, it becomes the new best known. A route
 537 is better if it has a higher score, or given the same score has a lower cost.
 538 Theoretically, this algorithm would compute the optimal solution if we run it
 539 without any additional constraints. However, as OP is an NP-hard problem, we
 540 want to keep the runtime and memory usage manageable, so we employ it as an

541 approximation algorithm. To keep memory under control it is possible to limit
542 the deque length: if, after inserting a new representative, the deque is longer
543 than the limit, the minimum element (the one with the lowest potential score)
544 is dropped.

545 As already mentioned, we only put a representative into the deque if its
546 potential score pot is higher than the best known route score, as such a partial
547 solution cannot possibly be expanded to an optimal solution. However, it is
548 also possible to prune more aggressively by ignoring representatives whose pot
549 is smaller than the best solution found so far multiplied by a *cut factor* t . In
550 other words, $best$ being the score of the best known solution, a representative
551 is only enqueued if $pot > t \cdot best$. In this way we may lose the optimal solution,
552 but since $extra$ is an overestimation we can compute a provable bound (see [5]
553 for a formal proof).

554 Limiting the length of the queue, we have to keep track of the highest po-
555 tential we ignored, let it be max ; then the optimal solution cannot be better
556 than $max/best$ times the score of the returned solution. When using a cutting
557 factor, t is the worst ratio between the scores of the optimal solution and the
558 returned one.

559 The approach is an anytime algorithm, since it continuously generates and
560 improves solutions. Thus, it is possible to keep runtime under control by re-
561 turning the current best known solution after reaching a timeout.

562 While this is an efficient approximation algorithm with provable bounds,
563 so far Bolzoni et al. assume that all the scores of the POIs are given as an
564 input parameter [6, 5]. Determining and custom-tailoring the scores of POIs
565 for individual users was left as an open problem. Here we solve this problem
566 by combining an efficient algorithm for solving the OP with categories [5] with
567 effective social sensing techniques, allowing us to automatically extract user
568 preferences and map them to POIs scores. More formally, given a set U of users
569 and a set P of POIs, the score of POI $p_j \in P$ for user $u_i \in U$ is calculated as
570 follows:

$$score(u_i, p_j) = SIM_{ij} \quad (24)$$

571 **5. Experimental Evaluation**

572 We implemented the overall pipeline described in Figure 2 and evaluated the
573 *accuracy* of the user interests in the POIs, the *user satisfaction* with reference
574 to the suggested paths, and *run-time performance* of both the *social sensing*
575 approach and the *orienteering* algorithm. In the following we provide some
576 details on the experimental setting before discussing the results.

577 *5.1. Environment*

578 We ran the experiments on a machine with a 2.30GHz Intel Core i5 processor
579 (3Mb of cache) with 8 GBytes of RAM running Archlinux, kernel version 4.9.6.

580 *5.2. Data Sets*

581 In order to evaluate the viability of our approach, we needed to construct
582 data sets of both *user tweets* and *POI descriptions*, which we describe in the
583 following.

584 *5.2.1. OSN Users*

585 Regarding the collected data set of OSN users, we downloaded the tweets of
586 549 users (on average 3200 tweets each) who are followers of the Municipality of
587 Bolzano Twitter account³. Then, we randomly selected 20 users from the users
588 who posted at least 100 tweets and who posted at least 60% of their tweets in
589 Italian.

590 All the user tweets were then preprocessed as described in Section 4.2, before
591 applying the schemes described in Section 4.

³These users have not put restrictions in their privacy settings, allowing us to do so.

592 *5.2.2. POI Descriptions*

593 In order to compare twitter accounts with POIs, we need to have a textual
594 representation of POIs. This dataset was constructed manually by searching
595 for textual descriptions of POIs in Bolzano that are interesting to a tourist.
596 The POI documents were then processed according to the schemes in Section 4.
597 Table 2 lists the set of considered POIs.

598 *5.3. Results*

599 First, we evaluate the *accuracy* of the social sensing approach (*RANK* ma-
600 trix) measured against a ground truth provided by human annotators. The
601 annotators were informed about a subset of the tweets of a subset of users.
602 They were asked to define a *top-k list* of POIs for each selected user based on
603 their tweets. Second, we measured the *user satisfaction* of each selected Twitter
604 user whom we suggested several tourist paths. Eventually, we also measured
605 the *processing time* of both the social sensing approaches and the *orienteering*
606 *algorithm*, illustrating how well this processing scales by using different numbers
607 of tweets and various time budgets.

608 *5.3.1. Evaluating Accuracy*

609 Given the set of users $U = \{u_i\}_{i \in [1, m]}$ (where $m \in \mathbb{N}$ is the total number
610 of users), the set of POIs $P = \{p_j\}_{j \in [1, n]}$ (where $n \in \mathbb{N}$ is the total number
611 of POIs), and the set of tweets $T = \{\lambda_i \in \Lambda \mid \lambda_i.s = "e"\}$ (Table 1), let $U' =$
612 $\{u'_i\}_{i \in [1, n_u]} \subset U$, $n_u \in \mathbb{N}$, $n_u < m$, $P' = \{p'_j\}_{j \in [1, n_p]} \subset P$, $n_p \in \mathbb{N}$, $n_p < n$, and
613 $T' = \{t'_l\}_{l \in [1, n_t]} \subset T$, $n_t \in \mathbb{N}$, $n < |T|$ be the number of users, the number of
614 POIs, and the number of tweets per user that n_{ann} annotators take into account
615 while performing the rankings, respectively.

616 Additionally, let k be the number of POIs for which both the system and
617 the annotators compute the rankings. More specifically, for each user $u \in U'$,
618 let $R_u^a = \{r_i^a\}_{i \in [1, k]}$ denote the ranks computed by the algorithm, which are the
619 POI IDs of the u -th row of the *RANK* matrix (R_{u*}); on the other hand, let

Table 2: Set P' of POIs

ID	Name of POIs
2000	South Tyrol Museum of Archaeology - Bolzano
2001	Bolzano cathedral - Bolzano
2002	Earth Pyramids Renon - Soprabolzano
2003	Neptune's Fountain - Bolzano
2004	Argentieri Road - Bolzano
2005	Chiesa dei Domenicani - Bolzano
2006	Messner Mountain Museum Firmian - Bolzano
2007	Franciscan Friary - Bolzano
2008	Museo di scienze naturali dell'Alto Adige - Bolzano
2009	Museion Museum of Modern and Contemporary Art - Bolzano
2010	Arcades - Bolzano
2011	Piazza delle Erbe - Bolzano
2012	Bolzano Station - Bolzano
3000	Renon cableway - Bolzano
3001	Runkelstein Castle - Bolzano
3002	Twenty Shopping Mall - Bolzano
3003	Victory Monument - Bolzano
3004	Maretsch Castle - Bolzano
3005	Emberfly Art Gallery - Bolzano
3006	Teatro comunale bolzano - Bolzano
3007	Palazzo Menz - Bolzano
3008	Antico Municipio - Bolzano
3009	Teutonic Order Church - Bolzano
3010	Monument to King Laurin - Bolzano
3011	Bottai Road - Bolzano
3012	Civic Museum - Bolzano

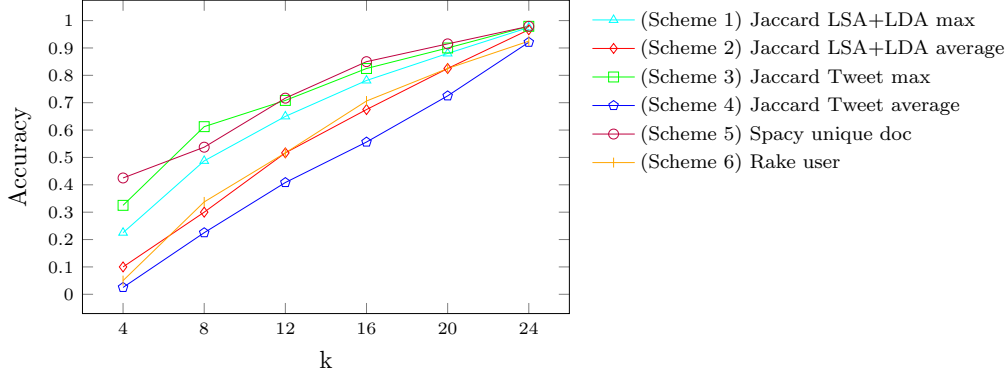


Figure 7: Accuracy results for our six schemes

620 $R_u^h = \{r_j^h\}_{j \in [1, k]}$ denote the ranks extracted by a human annotator. Then, for
 621 *each user and annotator*, the accuracy A is computed as follows:

$$A = \frac{|R_u^a \cap R_u^h|}{k}, \forall u \in U' \quad (25)$$

622 Clearly, A coincides with the classical definitions of both *precision* and *recall*,
 623 since both the system and the annotators are asked for *top-k lists* of POIs.

624 More specifically, in this use case we set $n_u = 20$, $n_p = 26$, $n_t = 3200$, and
 625 $n_{ann} = 3$. Table 2 shows the set P' of considered POIs.

626 Figure 7 shows the behavior of the average accuracy when varying k . More
 627 specifically, we reported the accuracy achieved by the six schemes which are
 628 *Jaccard LSA+LDA max*, *Jaccard LSA+LDA average*, *Jaccard Tweet max*, *Jac-*
 629 *card Tweet average*, *spaCy*, and *RAKE*. From Figure 7, we can infer that all
 630 the listed approaches have a similar linear increasing trend; *spaCy* and *Jaccard*
 631 *Tweet max* achieve better accuracy than the rest of the schemes. When k is
 632 relatively small ($k = 4$), *spaCy* outperforms *Jaccard Tweet max*. For $k = 8$,
 633 *Jaccard Tweet max* shows the best accuracy. Further increasing k makes the
 634 accuracy scores of *spaCy* and *Jaccard Tweet max* converge and nearly overlap.
 635 Between $k = 12$ and $k = 20$, there is no significant change in terms of the linear
 636 tendency and the relative positioning of the different schemes. As we will see,

637 the main advantage of *Jaccard Tweet* is that it can be significantly faster than
638 *spaCy*.

639 5.3.2. Evaluating User Satisfaction and Route Recall

640 In order to further evaluate the results, we have conducted a user study to
641 evaluate the *user satisfaction* and the *route recall* for the suggested routes. This
642 user study intends to evaluate (1) if the component of social sensing will increase
643 user satisfaction, (2) which of the six proposed score schemes leads to a higher
644 user satisfaction, and (3) the percentage of POIs in each route proposed to the
645 user that are actually listed in his/her top-k POI collection (i.e., computing the
646 *route recall*). Therefore, we derived nine customized routes for each user. Six of
647 them are based on the six score schemes, one is using fixed scores for each POI,
648 and, as baselines, the second to last one assigns scores randomly to the POIs,
649 whereas the last one just picks from the user ranking the top-k POIs which are
650 visitable within the time constraints.

651 We used active twitter accounts for sensing the user preference and POIs
652 from the city of Bolzano in Italy. As Bolzano is a touristic city and most POIs
653 are within walking distance, a user does not need to concern themselves with
654 any transportation issues. Also, the users in our study were not aware of the
655 different scoring schemes and also did not know that we analyzed their tweets.

656 We presented the following scenario to the user: assume you would like to
657 visit Bolzano and every POI there is new to you. Given a fixed total visiting
658 time of three hours and twenty minutes, you are not able to visit all the POIs.
659 We prepared a total of **nine** suggested routes. As an example, Table 4 shows
660 the routes proposed to the user *X*. Each route contains a set of POIs with
661 an orienteering route such as shown in Figure 8. The transportation means is
662 walking. After reviewing all the proposed routes and related POIs, we asked
663 users to rank the nine proposed routes based on their own preferences, thus
664 allowing the computation of the *user satisfaction*.

665 The generated routes were computed using the reference implementation of
666 the best first search algorithm in Section 4.6. Each pair of POI distances is

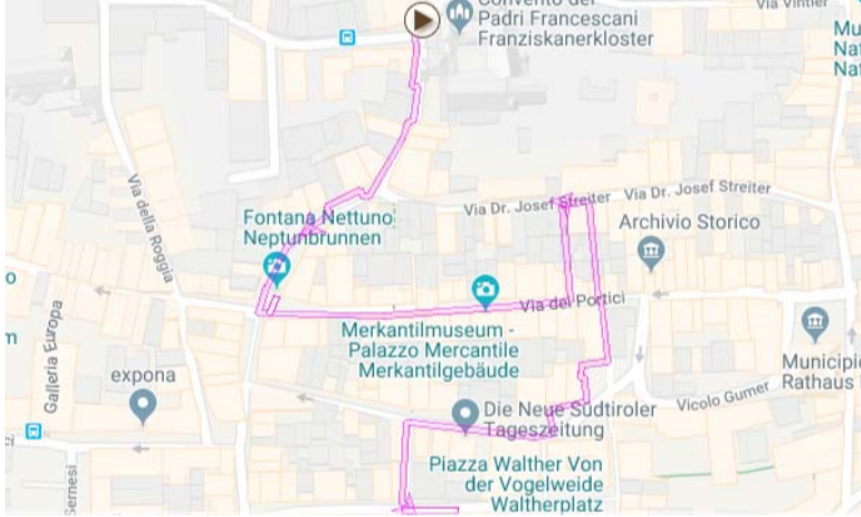


Figure 8: Route presented to the user

667 computed as the walking time distance according to Open Street Map and the
 668 POI scores are provided by social sensing. That leaves the budget time, starting,
 669 and ending points. We executed experiments for tours from 30 minutes up to
 670 200 minutes (3 hours and 20 minutes), here we focus on the results for the tours
 671 lasting 200 minutes. Furthermore, we assumed that to actually collect the score
 672 of a POI, a tourist needs to stay for 10 minutes at a POI. The tour starts and
 673 ends near the Franciscan Friary, which is a well-known and fairly central POI
 674 in Bolzano.

675 The 10 minutes visiting time is meant to be a bare minimum needed to
 676 actually get something out of a POI, but still long enough to make an impact
 677 on the budget time. The algorithm ensures that the destination is reached in
 678 less than the budget time, but in reality we expect the tours to last longer, as
 679 people will take pauses and stay in POIs longer than 10 minutes, yet even a 200
 680 minute tour is doable in an afternoon.

681 Thus, in order to compute the *route recall*, we first asked each user u to
 682 explicitly provide a ranked list of POIs $R_u^{gt} = \{r_i^{gt}\}_{i \in [1, n]}$, where gt stands for

683 *ground truth*. Then, assuming $P_u^a = \{p_i^a\}_{i \in [1, k]}$ to be the route (path) computed
 684 by the algorithm a for user u , and passing through the k -sequence of POIs
 685 $\{p_i^a\}_{i=1}^k$, the route recall for user u and path P_u^a is computed as follows:

$$R_{P_u^a} = \frac{|R_u^{gt} \cap P_u^a|}{k} \quad (26)$$

686 Thus, the *route recall* denotes the percentage of POIs in each route of length
 687 k matching the user top- k list of POIs.

688 For instance, Table 4 shows the routes proposed to the user X as well as the
 689 related *route recall*, that we computed by comparing the routes against the POI
 690 ranking list of user X shown in Table 3.

691 As a result, we can infer that for the user X , the route which achieves the
 692 highest recall is the (*Scheme 5*) *Spacy unique doc* (63.64%). In this case study,
 693 also the *random path* accidentally achieves a good accuracy, but this is clearly
 694 not valid in the general, average, case.

695 Additionally, the last row of Table 4 shows a path which just covers the
 696 top- k POIs of the list R_u^{gt} of the user u , where k is the maximum number of
 697 POIs visitable within the time constraints. Clearly, its route recall is 100%,
 698 but evidently its number of visited POIs is much lower with respect to the ones
 699 passed through by the other routes, since the path generation is not optimized
 700 in this case.

701 However, considering all the users, the results of the user study is that five
 702 out of six routes with social sensing are ranked higher than the three routes
 703 not exploiting social sensing, and also achieve a higher *route recall* (excluding
 704 the *Ideal Route*). This means, that adding social sensing to trip planning tends
 705 to have a positive impact on both user satisfaction and route recall. The five
 706 higher ranked routes are using the following schemes: spaCy unique doc, Jaccard
 707 LSA+LDA max, RAKE, Jaccard Tweet average, and Jaccard Tweet Max. The
 708 top three ranked routes are the variants spaCy unique doc (scheme 5), Jaccard
 709 LSA+LDA max (scheme 1), and RAKE (scheme 6), in that order. This is
 710 a strong indication that these three schemes are the top contenders when it

Table 3: POI Ranking List of User X

Rank	ID	Name of POIs
1	2011	Piazza delle Erbe - Bolzano
2	3000	Renon cableway - Bolzano
3	2010	Arcades - Bolzano
4	2007	Franciscan Friary - Bolzano
5	2001	Bolzano cathedral - Bolzano
6	3001	Runkelstein Castle - Bolzano
7	3004	Maretsch Castle - Bolzano
8	2005	Chiesa dei Domenicani - Bolzano
9	2000	South Tyrol Museum of Archaeology - Bolzano
10	2009	Museion Museum of Modern and Contemporary Art - Bolzano
11	2002	Earth Pyramids Renon - Soprabolzano
12	3008	Antico Municipio - Bolzano
13	3006	Teatro comunale bolzano - Bolzano
14	3011	Bottai Road - Bolzano
15	3007	Palazzo Menz - Bolzano
16	2004	Argentieri Road - Bolzano
17	3005	Emberfly Art Gallery - Bolzano
18	2003	Neptune's Fountain - Bolzano
19	2008	Museo di scienze naturali dell'Alto Adige - Bolzano
20	3012	Civic Museum - Bolzano
21	3010	Monument to King Laurin - Bolzano
22	3009	Teutonic Order Church - Bolzano
23	2006	Messner Mountain Museum Firmian - Bolzano
24	3003	Victory Monument - Bolzano
25	2012	Bolzano Station - Bolzano
26	3002	Twenty Shopping Mall - Bolzano

Table 4: The nine routes suggested to the user X

Route	POIs with Start and End 2007	Route Recall
(Scheme 1) Jaccard LSA+LDA max	⟨2008, 3000, 2012, 3007, 2005, 2009, 3012, 3003, 2000⟩	50%
(Scheme 2) Jaccard LSA+LDA average	⟨2010, 3007, 2009, 3005, 3003, 3012, 2011, 2005, 2012, 3009⟩	45.45%
(Scheme 3) Jaccard Tweet max	⟨2010, 3007, 3006, 2009, 3005, 3003, 3012, 2011, 2005, 2008⟩	45, 45%
(Scheme 4) Jaccard Tweet average	⟨2003, 3012, 2000, 2009, 2010, 2004, 2011, 2005, 2012, 3010⟩	54.54%
(Scheme 5) Spacy unique doc	⟨2011, 2005, 3006, 2001, 3007, 3008, 2010, 2000, 3012, 2009⟩	63.64%
(Scheme 6) Rake user	⟨2011, 2003, 2005, 2001, 3007, 3008, 2000, 3005, 3012⟩	50%
Fixed 500	⟨2011, 2003, 2010, 3008, 2004, 3007, 2005, 2001, 2008, 3011⟩	45.45%
Random	⟨2011, 2000, 3012, 2005, 2001, 3007, 2003, 3008, 3010, 3000⟩	54.54%
Ideal Route (Baseline)	⟨2011, 3000, 2010, 2007, 2001, 3001⟩	100%

711 comes to increasing the user satisfaction. The only social sensing scheme that is
712 not able to outperform the Fixed 500, Random, and Ideal Route approaches is
713 Jaccard LSA+LDA average. One of the reasons may be that the average of all
714 the Jaccard values introduces more topics that are not really favored by users.

715 We also conducted an interview with users to explore the reasons why cer-
716 tain routes are preferred. We found that if a POI that is strongly preferred by
717 a user is missing in a proposed tour, the user will rank the route lower. Also,
718 if some non-preferred POIs are on the path to interesting POIs, the user will
719 also consider to rank that route lower. As this is a walking tour, the environ-

720 ment plays a role as well: walking along a river or in a pedestrian zone usually
 721 increases the user satisfaction. Additionally, the Ideal Routes are ranked lower
 722 than the routes computed with the support of social sensing, since users prefer
 723 to visit as many POIs as possible within the time constraints and, even more
 724 important, do not want to follow paths which are not optimized in terms of total
 725 traveling time. Finally, we also report one of the evaluations of the approach
 726 in the field by one volunteer user who tested our framework for seven days in
 727 Bolzano. The feedback we received from this user was positive, which promises
 728 to make our approach applicable for real-world use.

729 5.3.3. Evaluating Performance

730 In this section, we evaluate the performance of the framework in terms of
 731 running times used by both the *social sensing* approach and the *orienteering*
 732 algorithm. Figure 9 shows how the six social sensing approaches scale when
 733 varying the number of analyzed tweets (each of the values shown on the x-axis
 734 of Figure 9 is the number of tweets of a specific user). In principle all of them
 735 achieve reasonable running times (times are in the order of seconds even when
 736 processing more than 2000 tweets), although the approaches shown in schemes
 737 1 and 2 are much slower than all the others. The fastest one is definitely the
 738 approach based on *spaCy* (Scheme 5), which can be even used to monitor and
 739 process large numbers of Tweets in real time.

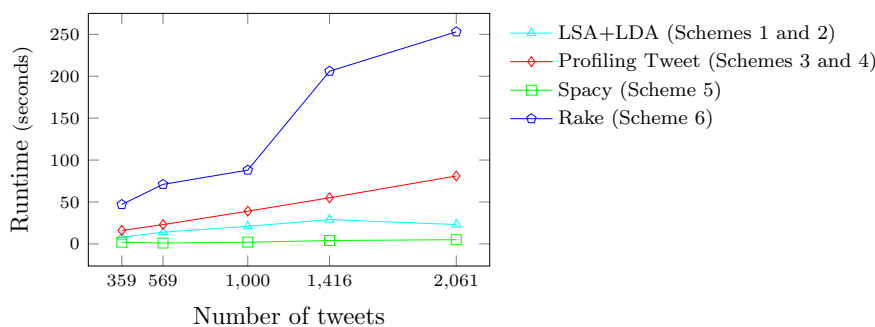


Figure 9: Run Time Performance - Social Sensing

740 Additionally, in the following we discuss the performance of the *orienteer-*
741 *ing algorithm* in terms of running time. We did not apply aggressive pruning
742 in the form of a cut factor, but we limited the run time to one minute. Fig-
743 ure 10(a) shows how the run time of the algorithm scales for different time
744 budgets, whereas Figure 11 refers to the run time of the orienteering algorithm
745 with that of the social sensing approaches. Even though we limited the run
746 time to one minute, the orienteering algorithm can go beyond that, since it only
747 stops at the end of a complete expansion phase.

748 It is also interesting to note that the run time does not grow monotonically
749 with the budget time. The reason for this is that the POIs are not distributed
750 evenly throughout the city, which means that for different budgets different sets
751 of POIs can be reached, making it more (or less) difficult to maximize the overall
752 score of a route.

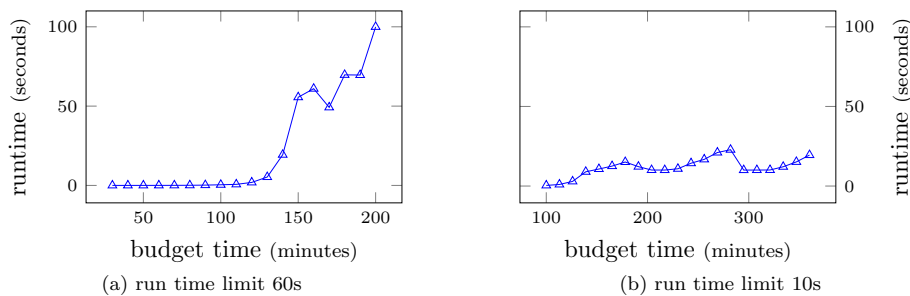


Figure 10: Route generation time, increasing the time budget.

753 If a fast response time is important, we can also lower the run time limit
754 of the orienteering algorithm. Figure 10(b) shows how the algorithm scales for
755 different budgets with a run time limit of ten seconds (up to routes with a length
756 of 360 minutes). Reducing the response time to ten seconds only has a minor
757 impact on the quality of the generated routes in terms of the overall score. On
758 average, the scores of the routes produced by a ten-second execution were only
759 about 1% lower than those generated by running the orienteering algorithm for
760 one minute.

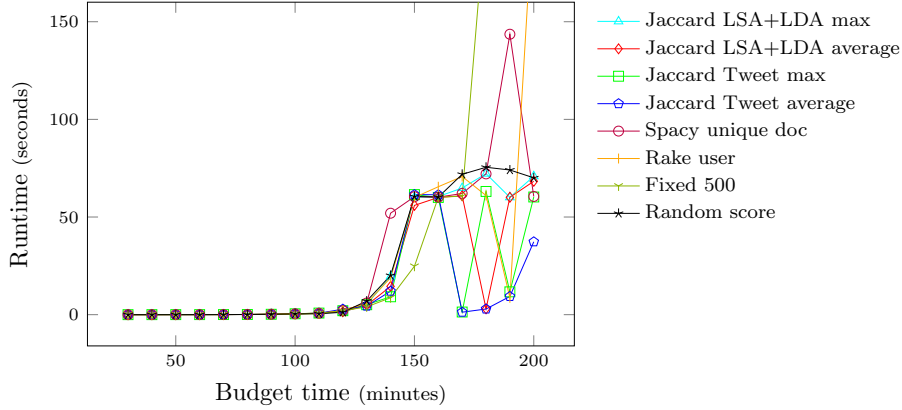


Figure 11: Runtime of each score

761 5.3.4. Final Remarks

762 Overall, we were able to effectively integrate all of the six proposed ap-
763 proaches for social sensing with the algorithm for solving the orienteering prob-
764 lem with category constraints, and achieve good performance in terms of *accu-*
765 *racy*, *user satisfaction*, and *efficiency*; in fact, five out of the six approaches are
766 ranked by the users better than the orienteering algorithm not exploiting social
767 sensing. However, it turns out that the approach which performs best is the one
768 based on *spaCy* (scheme 5).

769 Furthermore, in personalized trip planning, users may expect to explore new
770 or surprising travel experiences that usually cannot be captured by stable user
771 preferences [42]. In this paper, we only consider and evaluate route quality
772 as the alignment between the personalized POIs and user preferences. Other
773 evaluation factors such as serendipity or coverage [43] of POIs are out of the
774 scope of this paper. Since the main goal of the tourist trip planning is to
775 recommend personalized routes for visiting highly preferred POIs, we provided
776 each user with eight possible routes to choose from (six of which derived the
777 user preferences by social sensing).

778 6. Conclusions

779 One widely used preference elicitation technique is to explicitly ask if a
780 user likes or dislikes a certain item that was previously experienced by the
781 user. In recommender systems, the preference for an item is usually represented
782 by a rating for the item [44]. Although from the user’s perspective explicit
783 preference elicitation can be seen as a redundant procedure with considerable
784 effort on the side of the user, it is still employed due to its effectiveness in
785 many application domains [45]. However, there is always the risk of users losing
786 interest in a system because of the effort they have to invest to communicate
787 their preferences. Thus, we advocate the use of implicit preference elicitation
788 whenever possible (i.e., whenever the quality of a service does not suffer) to
789 keep the overhead as small as possible. The advent of social sensing techniques
790 provides us with powerful tools for automating preference elicitation.

791 More concretely, we propose an innovative method for customizing the so-
792 lutions to orienteering problems by applying six different schemes for social
793 sensing. In tourist trip planning, we identify user interests by sensing their
794 Twitter profile and are able to create an ordered list of POIs corresponding to
795 these interests. In this way, an orienteering algorithm can be custom-tailored to
796 a user by considering the semantic closeness between POIs and a user’s social
797 contents.

798 We implemented a prototype and evaluated its accuracy, user satisfaction
799 levels, and computational performance. The evaluation has shown the effective-
800 ness of exploiting social sensing for improving tourist trip planning. In partic-
801 ular, the approach that has shown the best user satisfaction is the one based
802 on the *spaCy* tool. Thus, the evaluation demonstrates that it is very efficient
803 and promising to integrate social sensing with orienteering-based tourist trip
804 planning, making it applicable to real-world use cases.

805 In the future we want to build on this work by making it scale to even
806 larger numbers of users and POIs, and we want to extend our approach to other
807 languages and additional social networks, such as Facebook or Instagram. We

808 do not see switching to another Indo-European language as an issue, since, for
809 instance, the fastText framework [32] is available for 157 different languages.

810 References

- 811 [1] S. Liu, L. Wang, A self-adaptive point-of-interest recommendation algo-
812 rithm based on a multi-order markov model, *Future Generation Computer*
813 *Systems* 89 (2018) 506–514. doi:10.1016/j.future.2018.07.008.
- 814 [2] M. Eirinaki, J. Gao, I. Varlamis, K. Tserpes, Recommender systems for
815 large-scale social networks: A review of challenges and solutions, *Fu-*
816 *ture Generation Computer Systems* 78 (2018) 413–418. doi:10.1016/j.
817 *future.2017.09.015*.
- 818 [3] S. Milovanovi, Z. Bogdanovi, A. Labus, D. Bara, M. Despotovi-Zraki,
819 An approach to identify user preferences based on social network anal-
820 ysis, *Future Generation Computer Systems* 93 (2018) 121 – 129. doi:
821 10.1016/j.future.2018.10.028.
- 822 [4] S. Rose, D. Engel, N. Cramer, W. Cowley, Automatic keyword extraction
823 from individual documents, in: M. W. Berry, J. Kogan (Eds.), *Text Mining.*
824 *Applications and Theory*, John Wiley and Sons, Ltd, 2010, pp. 1–20. doi:
825 10.1002/9780470689646.ch1.
- 826 [5] P. Bolzoni, S. Helmer, Hybrid best-first greedy search for orienteering with
827 category constraints, 2017, pp. 24–42. doi:10.1007/978-3-319-64367-0_
828 2.
- 829 [6] P. Bolzoni, F. Persia, S. Helmer, Itinerary planning with category con-
830 straints using a probabilistic approach, in: D. Benslimane, E. Damiani,
831 W. I. Grosky, A. Hameurlain, A. Sheth, R. R. Wagner (Eds.), *Database and*
832 *Expert Systems Applications*, Springer International Publishing, Cham,
833 2017, pp. 363–377. doi:10.1007/978-3-319-64471-4_29.

- 834 [7] T. Tsiligirides, Heuristic methods applied to orienteering, *The Journal of*
835 *the Operational Research Society* 35 (9) (1984) 797–809. doi:10.1057/
836 *jors.1984.162*.
- 837 [8] E. H. Lu, C. Lin, V. S. Tseng, Trip-mine: An efficient trip planning ap-
838 *proach with travel time constraints*, in: *2011 IEEE 12th International*
839 *Conference on Mobile Data Management*, Vol. 1, 2011, pp. 152–161.
840 doi:10.1109/MDM.2011.13.
- 841 [9] G. Righini, M. Salani, Decremental state space relaxation strategies and ini-
842 *tialization heuristics for solving the orienteering problem with time windows*
843 *with dynamic programming*, *Computers & Operations Research* 36 (4)
844 (2009) 1191 – 1203. doi:https://doi.org/10.1016/j.cor.2008.01.003.
- 845 [10] C. Keller, Algorithms to solve the orienteering problem: A comparison,
846 *European Journal of Operational Research* 41 (2) (1989) 224 – 231. doi:
847 [https://doi.org/10.1016/0377-2217\(89\)90388-3](https://doi.org/10.1016/0377-2217(89)90388-3).
- 848 [11] A. Blum, S. Chawla, D. R. Karger, T. Lane, A. Meyerson, M. Minkoff,
849 *Approximation algorithms for orienteering and discounted-reward tsp*, in:
850 *44th Annual IEEE Symposium on Foundations of Computer Science*, 2003.
851 *Proceedings.*, 2003, pp. 46–55. doi:10.1109/SFCS.2003.1238180.
- 852 [12] C. Chekuri, N. Korula, M. Pál, Improved algorithms for orienteering and
853 *related problems*, *ACM Trans. Algorithms* 8 (3) (2012) 23:1–23:27. doi:
854 [10.1145/2229163.2229167](https://doi.org/10.1145/2229163.2229167).
- 855 [13] Z. Sevkli, F. E. Sevilgen, Variable neighborhood search for the orienteering
856 *problem*, in: A. Levi, E. Savaş, H. Yenigün, S. Balcısoy, Y. Saygın (Eds.),
857 *Computer and Information Sciences – ISCIS 2006*, Springer Berlin Heidel-
858 *berg*, Berlin, Heidelberg, 2006, pp. 134–143. doi:10.1007/11902140_16.
- 859 [14] M. F. Tasgetiren, A. E. Smith, A genetic algorithm for the orienteering
860 *problem*, in: *Proceedings of the 2000 Congress on Evolutionary Com-*

- 861 putation. CEC00 (Cat. No.00TH8512), Vol. 2, 2000, pp. 910–915 vol.2.
862 doi:10.1109/CEC.2000.870739.
- 863 [15] Q. Wang, X. Sun, B. L. Golden, J. Jia, Using artificial neural networks
864 to solve the orienteering problem, *Annals of Operations Research* 61 (1)
865 (1995) 111–120. doi:10.1007/BF02098284.
- 866 [16] P. Bolzoni, S. Helmer, K. Wellenzohn, J. Gamper, P. Andritsos, Efficient
867 itinerary planning with category constraints, in: *23rd ACM SIGSPATIAL*
868 *Int. Conf. on Advances in Geographic Information Systems (ACM SIGSPA-*
869 *TIAL'14)*, 2014, pp. 203–212. doi:10.1145/2666310.2666411.
- 870 [17] W. Huang, I. Weber, S. Vieweg, Inferring nationalities of twitter users and
871 studying inter-national linking, in: *Proceedings of the 25th ACM Confer-*
872 *ence on Hypertext and Social Media, HT '14*, ACM, New York, NY, USA,
873 2014, pp. 237–242. doi:10.1145/2631775.2631825.
- 874 [18] J. Bentwood, *Distributed influence: Quantifying the impact of social media*,
875 Edelman White Paper (2007).
- 876 [19] R. Tinati, L. Carr, W. Hall, J. Bentwood, Identifying communicator roles
877 in twitter, in: *Proceedings of the 21st International Conference on World*
878 *Wide Web, WWW '12 Companion*, ACM, New York, NY, USA, 2012, pp.
879 1161–1168. doi:10.1145/2187980.2188256.
- 880 [20] E. D'Avanzo, G. Pilato, M. Lytras, Using twitter sentiment and emo-
881 tions analysis of google trends for decisions making, *Program* 51 (3)
882 (2017) 322–350. arXiv:<https://doi.org/10.1108/PROG-02-2016-0015>,
883 doi:10.1108/PROG-02-2016-0015.
- 884 [21] S. Gianvecchio, M. Xie, Z. Wu, H. Wang, Measurement and classification of
885 humans and bots in internet chat, in: *Proceedings of the 17th Conference*
886 *on Security Symposium, SS'08*, USENIX Association, Berkeley, CA, USA,
887 2008, pp. 155–169.
888 URL <http://dl.acm.org/citation.cfm?id=1496711.1496722>

- 889 [22] B. Sarwar, G. Karypis, J. Konstan, J. Riedl, Item-based collaborative filter-
890 ing recommendation algorithms, in: Proceedings of the 10th International
891 Conference on World Wide Web, WWW '01, ACM, New York, NY, USA,
892 2001, pp. 285–295. doi:10.1145/371920.372071.
- 893 [23] T. K. Landauer, P. W. Foltz, D. Laham, An introduction to latent se-
894 mantic analysis, *Discourse processes* 25 (1998) 259 – 284. doi:10.1080/
895 01638539809545028.
- 896 [24] T. K. Landauer, S. T. Dumais, A solution to plato’s problem: The latent
897 semantic analysis theory of acquisition, induction, and representation of
898 knowledge, *Psychological review* 104(2) (1997) 211 – 223. doi:10.1037/
899 /0033-295X.104.2.211.
- 900 [25] F. Agostaro, A. Augello, G. Pilato, G. Vassallo, S. Gaglio, A conversational
901 agent based on a conceptual interpretation of a data driven semantic space,
902 *Lecture Notes in Artificial Intelligence* 3673(2) (2005) 381–392. doi:10.
903 1007/11558590_39.
- 904 [26] S. Santilli, L. Nota, G. Pilato, The use of latent semantic analysis in the
905 positive psychology: A comparison with twitter posts, *Semantic Computing*
906 (ICSC), 2017 IEEE 11th International Conference on IEEE (2017) (2017)
907 494498. doi:10.1109/ICSC.2017.60.
- 908 [27] D. Blei, A. Ng, M. Jordan, Latent dirichlet allocation, *Journal of Machine*
909 *Learning Research* 3 (January 2003) 9931022.
- 910 [28] T. Hofmann, Probabilistic latent semantic indexing, in: Proceedings of
911 the 22Nd Annual International ACM SIGIR Conference on Research and
912 Development in Information Retrieval, SIGIR '99, ACM, New York, NY,
913 USA, 1999, pp. 50–57. doi:10.1145/312624.312649.
- 914 [29] D. M. Blei, Probabilistic topic models, *Commun. ACM* 55 (4) (2012) 77–84.
- 915 [30] W. Darling, A theoretical and practical implementation tutorial on topic
916 modeling and gibbs sampling, Proceedings of the 49th annual meeting of

- 917 the association for computational linguistics: Human language technologies
918 (Dec 2011) 642–647.
- 919 [31] Y. Whye Teh, D. Newman, M. Welling, A collapsed variational bayesian
920 inference algorithm for latent dirichlet allocation, Proceedings of NIPS 6
921 (2006) 1378–1385.
- 922 [32] fasttext : Library for efficient text classification and representation learning
923 (2018).
924 URL <https://fasttext.cc/>
- 925 [33] spacy, the industrial-strength natural language processing (2018).
926 URL <https://spacy.io/>
- 927 [34] F. Persia, S. Helmer, S. Pugacs, G. Pilato, Social sensing for improving the
928 user experience in orienteering, in: 2019 IEEE 13th International Confer-
929 ence on Semantic Computing (ICSC), 2019. doi:10.1109/ICSC.2019.
930 8665498.
- 931 [35] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed rep-
932 resentations of words and phrases and their compositionality, in: Advances
933 in neural information processing systems, 2013, pp. 3111–3119.
- 934 [36] E. Grave, P. Bojanowski, P. Gupta, A. Joulin, T. Mikolov, Learning word
935 vectors for 157 languages, arXiv preprint arXiv:1802.06893.
- 936 [37] F. Persia, S. Helmer, A framework for high-level event detection in a social
937 network context via an extension of iseq1, in: 2018 IEEE 12th International
938 Conference on Semantic Computing (ICSC), 2018, pp. 140–147. doi:10.
939 1109/ICSC.2018.00028.
- 940 [38] F. Amato, A. Santo, V. Moscato, F. Persia, A. Picariello, Detecting
941 unexplained human behaviors in social networks, in: 2014 IEEE Inter-
942 national Conference on Semantic Computing, 2014, pp. 143–150. doi:
943 10.1109/ICSC.2014.21.

- 944 [39] F. Amato, A. Castiglione, A. D. Santo, V. Moscato, A. Picariello, F. Persia,
945 G. Sperl, Recognizing human behaviours in online social networks, *Com-
946 puters & Security* 74 (2018) 355 – 370. doi:[https://doi.org/10.1016/
947 j.cose.2017.06.002](https://doi.org/10.1016/j.cose.2017.06.002).
- 948 [40] G. Pilato, G. Vassallo, Tsvd as a statistical estimator in the latent semantic
949 analysis paradigm, *IEEE Transactions on Emerging Topics in Computing*
950 3 (2) (2015) 185–192. doi:[10.1109/TETC.2014.2385594](https://doi.org/10.1109/TETC.2014.2385594).
- 951 [41] A. Cuzzocrea, G. Pilato, Taxonomy-based detection of user emotions for
952 advanced artificial intelligent applications, *Hybrid Artificial Intelligent
953 Systems, Lecture Notes in Computer Science* 10870 (2018) 573 – 585.
954 doi:[10.1007/978-3-319-92639-1_48](https://doi.org/10.1007/978-3-319-92639-1_48).
- 955 [42] J. Xu, J. Chen, R. Zhou, J. Fang, C. Liu, On workflow aware location-
956 based service composition for personal trip planning, *Future Generation
957 Computer Systems* 98 (2019) 274 – 285. doi:[10.1016/j.future.2019.
958 03.010](https://doi.org/10.1016/j.future.2019.03.010).
- 959 [43] M. Ge, C. Delgado-Battenfeld, D. Jannach, Beyond accuracy: evaluating
960 recommender systems by coverage and serendipity, in: *Proceedings of the
961 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona,
962 Spain, September 26-30, 2010, 2010*, pp. 257–260. doi:[10.1145/1864708.
963 1864761](https://doi.org/10.1145/1864708.1864761).
- 964 [44] R. Logesh, V. Subramaniaswamy, V. Vijayakumar, X.-Z. Gao, V. Indra-
965 gandhi, A hybrid quantum-induced swarm intelligence clustering for the
966 urban trip recommendation in smart city, *Future Generation Computer
967 Systems* 83 (2018) 653 – 673. doi:[10.1016/j.future.2017.08.060](https://doi.org/10.1016/j.future.2017.08.060).
- 968 [45] A. Pommeranz, J. Broekens, P. Wiggers, W.-P. Brinkman, C. M. Jonker,
969 Designing interfaces for explicit preference elicitation: a user-centered in-
970 vestigation of preference representation and elicitation process, *User Mod-
971 eling and User-Adapted Interaction* 22 (4) (2012) 357–397. doi:[10.1007/
972 s11257-011-9116-6](https://doi.org/10.1007/s11257-011-9116-6).