

RESEARCH ARTICLE

A Place Recommendation Approach Using Word Embeddings in Conceptual Spaces

OMID R. ABBASI¹ AND ALI A. ALESHEIKH¹

Department of GIS, K. N. Toosi University of Technology, Tehran 19967, Iran

Corresponding author: Omid R. Abbasi (oabbasi@mail.kntu.ac.ir)

Many thanks to my friend in Switzerland for covering the publication costs.

ABSTRACT The way that computing systems digest geographic space is fundamentally different from people's understanding of space. In human discourse, a geographic space is referred to by a place name, and the reasoning about a place are based on its characteristics. This is in contrast with computing systems where geographical spaces are handled by the definition of coordinate systems. Hence, when recommending places, a recommendation method that leverages textual content, as a medium of communication among people, can be better understood. In this paper, we use elements of Natural Language Processing (NLP), such as Positive Point-wise Mutual Information (PPMI), Term Frequency - Inverse Document Frequency (TF-IDF), and Multi-Dimensional Scaling (MDS), to infer a conceptual space of the items of a place-based recommender system. By applying a Support Vector Machine (SVM) classifier on the resulting conceptual space, some meaningful directions are extracted. Shannon entropy is used as a measure to identify the directions that imply a valid geographic region. We apply the method on a dataset of advertisement descriptions of rental properties and a dataset of Persian Wikipedia articles. The results showed the proposed method is able to measure the similarity of items in the inferred conceptual space with 88% of accuracy. A comparison with BERT algorithm demonstrates the superiority of the proposed method over the baseline models.

INDEX TERMS Conceptual spaces, place, recommender systems, semantic similarity, textual content, word embeddings.

I. INTRODUCTION

The extraction of meaning from textual content is used in various research areas and applications, including information retrieval [1], sentiment analysis [2], event detection [3], and recommender systems [4], [5], [6]. Some approaches utilize Natural Language Processing techniques along with Machine Learning to deduce meaningful information [7], [8], [9]. Today, users can express their opinions about every topic on the social media, discuss about their common interests on the chat rooms, and advertise their own properties on an online market. Text is the main means of communication among people over the World Wide Web [10], [11]. Place-related information is concealed in plenty of these unstructured data. Places are the prevalent communication means of individuals when referring to space [12], [13]. People, unlike computing systems (e.g. Geographic Information Systems (GIS)), do not

use coordinates to refer to the objects in the space. They refer to those entities using a place name [14]. The focus of many studies has been on the extraction spatial relations and identifying the location of place names in texts [15]. However, only 0.2% of sentences in textual resources contain spatial descriptors [16]. Hence, it is vital for place-based recommenders to extract semantics that imply space.

In this paper, we argue that some information related to places can be deduced from non-spatial terms. For example, *well-off* is a place-based term when it is stated about the neighborhoods of a city. If a user is going to rent an apartment in a well-off neighborhood, an ideal recommender system should be able to suggest other properties in a similar neighborhood. Due to the rare occurrence of spatial references in texts, it is vital for place-based recommenders to extract semantics that imply space.

The aim of this study is to extract meaningful place-based concepts from texts and to suggest similar places, which is applicable in recommender systems. The extraction of

The associate editor coordinating the review of this manuscript and approving it for publication was Le Hoang Son¹.

place-based concepts is based on the non-spatial terms within texts. In general, the advantages of the proposed approach are two-fold:

- The proposed approach leverages components from the cognitive theory of conceptual spaces. Hence, the similarity of items is computed in a cognitive space, which is essential in a computing system focusing on human users.
- The proposed approach is fully unsupervised, making it suitable in scenarios where labeled data are limited or not accessible.

The contributions of this study are three-fold:

- The vast majority of studies conducted in this research area have focused on extracting spatial terms such as place names and spatial relations from texts. The proposed approach utilizes semantics of terms to induce spatial properties of textual items.
- We extend the original theory of conceptual spaces to generate a cognitive space suitable for computing similarities in recommenders. We also develop a method to measure the similarity of items (places). The proposed similarity measure is directional and is capable of measuring the dissimilarity of items.
- We use Shannon entropy as a measure to identify the terms that imply spatial connotations. This is necessary for the proposed method, because the extracted topics do not necessarily refer to a specific geographical region in the study area.

The paper is organized as follows. In Section II, some related work conducted on the intersection of places and recommenders are briefly introduced. Then, our proposed approach towards constructing the conceptual space from textual content is presented. In Section IV, the approach is applied on a real-world dataset and the results are presented and discussed. Finally, Section V concludes the paper and presents future insights.

II. RELATED WORK

The literature is quite rich in terms of the developments achieved in recommender algorithms and place-based research. Since this is a vast research area and includes a wide range of methods, this section is focused on those studies that handle textual content or utilize semantic similarity to provide place recommendations.

Traditionally, recommender systems gather their data from items and/or users. A great number of research have utilized data about places, and have computed similarities based on a variety of computational approaches such as meta-heuristics and artificial intelligence methods [19]. The studies have commonly employed either the location of users or the location of items, such as geo-tagged images [20], [21], as the core part of their methodology and recommend items ranked higher by near users to the current user. Context, as another dimension, has been integrated to recommender systems in many researches. Some studies have utilized ambient data,

such as weather information [22], traffic condition [23], temporal dimension [24], or even more fine-grained information such as user orientation [25], as a proxy for context, and some others have inferred the context from content. The latter includes analyzing image content [26] and text mining [27], [28], among many others. Capdevila et al. [29] have developed a recommender approach based on text mining techniques. They have crawled foursquare social network and gathered the user reviews on places. By applying sentiment and content analyses, they have proposed a hybrid recommender based on which similar pairs of users and items are identified. Amara and Subramanian [30] have tackled the cold start problem and constructed user profiles from analyzing textual content. Zhao et al. [31] have combined semantic, textual, and location information to train a deep learning model with the aim of identifying appropriate jobs for users. Tao et al. [32] have combined sentiment analysis and topic modeling techniques to enhance the recommendation of places of interest (POIs). Bafna et al. [33] have used ontologies to compute semantic similarities between terms in news. They have constructed a feature matrix where each row represents a feature vector for a document. Then, they have applied a clustering algorithm on the matrix based on semantic similarities extracted from WordNet. The use of ontologies to extract semantic relations has been extensively studied in other studies as well [34] and [35].

In some other studies, based on available information, word embeddings are generated. For instance, Vasile et al. [36] have proposed a method, called Meta-Prod2Vec, to compute item similarities using their metadata. Their approach utilizes user interactions with items and their attributes to train an ANN. They have applied their method on a dataset of music and concluded that their approach outperforms commonly used Prod2Vec method. Liu et al. [37] have extended the continuous bag of words model (CBoW) to capture the sentiment analysis and the domain to which each word belongs. They have applied their model on a dataset of Amazon user reviews, and achieved an improvement of accuracy by 2%. Faruqui, et al. [38] have proposed a method to improve the embedding representation by leveraging the relational knowledge (e.g. by knowledge graphs) extracted from ontologies. They have compared the results of their proposed method against the state-of-the-art approaches in semantic embeddings.

III. CONCEPTUAL SPACES

The theory of conceptual spaces [39] originates from cognitive sciences. It is an intermediate framework of knowledge representation, lying between symbolic and connectionist approaches [40]. It leverages geometrical representation of concepts in a high-dimensional space and measures the similarity between concepts or instances of concepts by a metric defined on the space. This theory is advantageous over using symbols as it decomposes the concepts into their semantic characteristics and, to some extent, addressing the

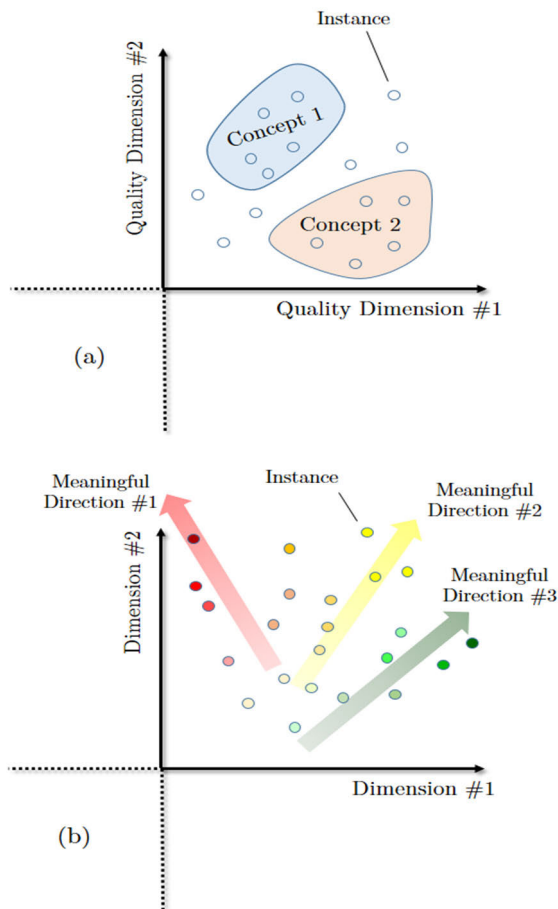


FIGURE 1. The schematic comparison of representing concepts in the original conceptual spaces theory and the proposed method.

symbol grounding problem. In addition, the structure and the interactions of the framework are more transparent than that of the connectionist models, enabling the interpretability of the mechanism of the framework. In this paper, the original theory is extended and some features over the conceptual space are developed. While the theory represents the concepts as convex regions in a high-dimensional space [41], we deal with semantic information as meaningful directions inside the space. In addition, the dimensions of the space in the original theory, called quality dimensions, are determined based on the decomposed qualities of the concepts [42], [43]. However, in most cases, no prior knowledge about the defining qualities of concepts is available. Some researchers have focused on data-driven approaches to induce the structure of the conceptual space [44], [45]. In current study, the space is constructed from the dataset without any prior knowledge about the place-based concepts. Fig. 1 schematically compares the differences between some aspects of the original theory (Fig. 1a) and that of our data-driven approach (Fig. 1b).

IV. PROPOSED METHOD

The goal of the proposed method is to generate a data-driven similarity space in which meaningful directions can be

identified. While the initial works on Gärdenfors' conceptual spaces mainly include explicitly defined quality dimensions, there have been some efforts recently to infer the conceptual space in a data-driven approach. In our method, we aim to construct a conceptual space resulting from the vector representation of a corpus. Fig. 2 illustrates the workflow of our methodology. The following subsections introduce the different parts of the workflow.

A. DATA COLLECTION

From a geographical point of view, places are spaces that are intertwined with meaning. Therefore, geographical coordinates alone cannot be the identity of a place. In place-based research, the meaning and characteristics of a place are usually extracted from its description. With the emergence and expansion of social networks over the web, a large amount of information related to places is available. The method presented in this research is based on the processing of textual information related to the places. Therefore, the first step is to collect textual documents related to the places under study. In order to evaluate the study, it is necessary for the datasets to contain the geographical coordinates of the place in addition to the textual information.

B. DATA CLEANING

As with all data-driven studies, it is necessary to check the validity of the dataset. In this step, the extracted textual information are examined to determine whether they meet the expectations of the research or not. Specifically, if a place has no description, it should be removed from the dataset. Depending on the use case, even if the description is very short, it can be left out of the dataset, because very short descriptions usually do not contain information that help to know more about that place. Moreover, there is a possibility that the geographic coordinates are entered incorrectly and are outside the study area.

C. DATA PREPROCESSING

In the pre-processing step, the documents are emptied of content that lacks semantic value related to the places. Typically, these contents include punctuations, alphanumeric characters, and stop words. Depending on the use case, some parts of speech can also be omitted. For example, in a specific case study such as real estate advertisements, the range of verbs used is limited and mostly does not refer to the characteristics of the place. On the other hand, nouns and adjectives can be very helpful in fathoming the characteristics of the place. After removing unnecessary terms, the documents are tokenized and stored as bags of words (BoW).

D. CONCEPTUAL SPACE CONSTRUCTION

The purpose of this step is to create a similarity space where similar documents are located at close distances to each other. In order to create such a space, it is first necessary to quantify the documents. Methods such as TF-IDF and PPMI have been used in various studies for this purpose. Given a set D of

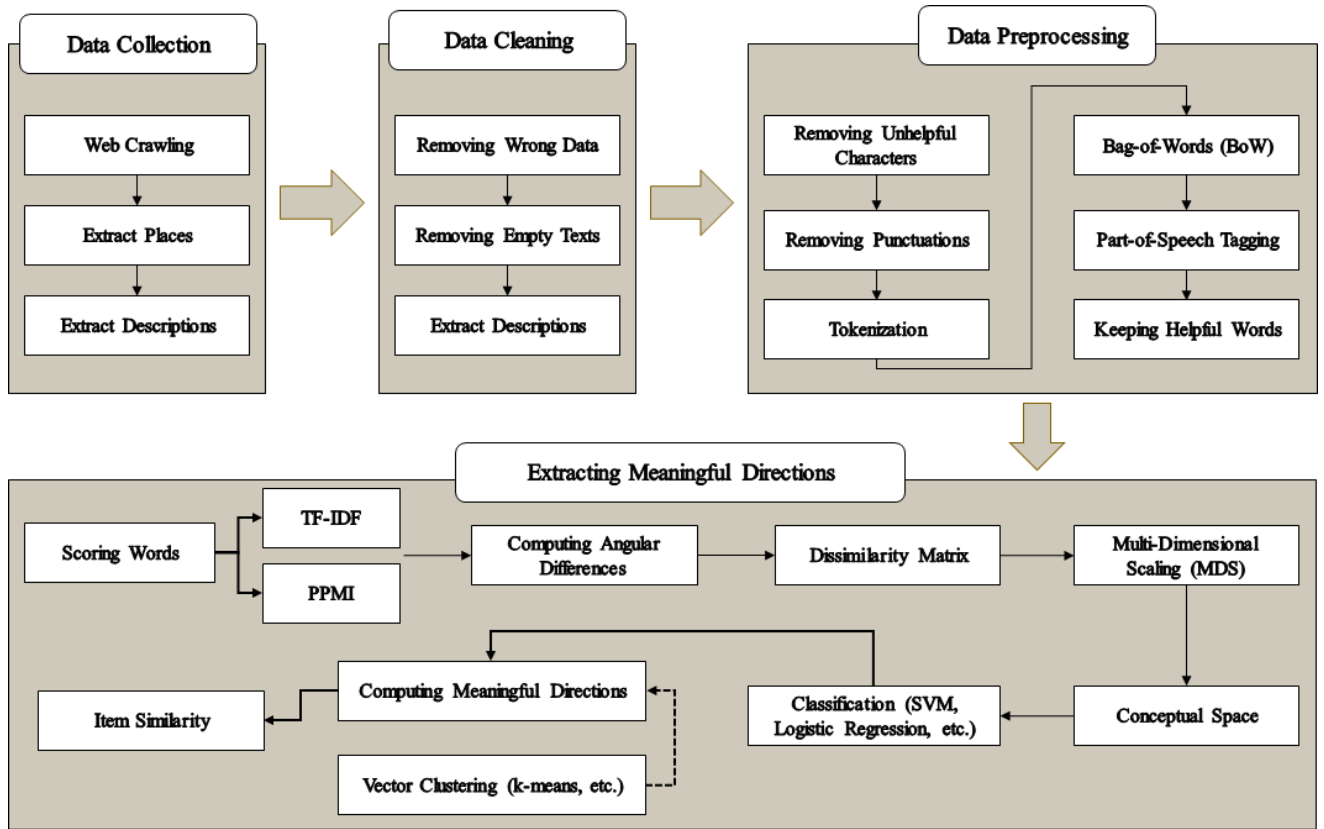


FIGURE 2. The workflow of the proposed approach towards finding meaningful directions in conceptual spaces.

documents, TF-IDF scores the word w in the document d as Eq. (1).

$$TF - IDF(w, d, D) = TF(w, d) \times IDF(w, D) \quad (1)$$

where

$$TF(w, d, D) = \frac{c(w, d)}{|d|} \quad (2)$$

$$IDF(w, D) = \log \left(\frac{|D|}{|\{d \in D | w \in d\}|} \right) \quad (3)$$

where $c(w, d)$ is the number of times term w occurs in document d . Also, $|d|$ and $|D|$ denote the number of words in document d and the number of documents in corpus D , respectively.

PPMI is similar to TF-IDF as it also scores the words in a document. Some researchers have shown that it works better than other scoring methods for the purpose of semantic similarity [46]. The index is calculated as Eq. (4).

$$PPMI(w, d) = \max \left(0, \log \left(\frac{p_{w,d}}{p_{w,*} \times p_{*,d}} \right) \right) \quad (4)$$

where

$$p_{w,d} = \frac{c(w, d)}{\sum_{w'} \sum_{d'} c(w', d')} \quad (5)$$

$$p_{w,*} = \sum_{d'} p_{w,d'} \quad (6)$$

$$p_{*,d} = \sum_{w'} p_{w',d} \quad (7)$$

By applying the scoring method, for each document d a vector v_d is formed. This vector contains the index values of all words in the corpus. By calculating the angle between these vectors, their similarity is obtained. Finally, by using the MDS method, the coordinates of the documents, as points in the conceptual space, are produced. The goal of this technique is to achieve the highest possible correlation between dissimilarity values and the distances among points. Unlike traditional conceptual spaces where dimensions have a pre-determined meaning, the interpretation of the dimensions of the output of MDS is not trivial [47], [48]. In addition, it is not even clear how many dimensions are needed for a suitable conceptual space. Therefore, it is a trial and error process and should be accomplished by different number of dimensions. A criterion for determining the optimum number of dimensions is maximizing the accuracy of the classification used in the next step.

E. IDENTIFYING MEANINGFUL DIRECTIONS

By applying MDS, those documents that are similar to each other are located closer to each other. In order to find meaningful directions, a classifier (e.g. Support Vector Machine (SVM) or Logistic Regression (LR) Classifier) is used to partition the space by a hyper-plane for each word. The classifier aims to distinguish those points (documents) containing a given word from points (documents) lacking the given word. Therefore, before each classification, the documents are labeled as whether they contain the given word or not. The hyper-plane should be placed such that those documents containing word w be grouped on one side of the plane. Then, the vector S_w perpendicular to this plane shows a meaningful direction. As this classification does not potentially yield a perfect accuracy, only those words for which the classification has a promising accuracy are taken into account. This is to ascertain the found vectors have strong meanings. However, even after constructing the hyper-plane only for those words with a high accuracy of classification, a high number of vectors with similar directions might exist in the conceptual space. As a similarity space, the vectors with close meanings are pointing towards a close direction. The vectors can be clustered to achieve more robust and meaningful directions within the constructed space.

V. EXPERIMENT AND DISCUSSION

A. IMPLEMENTATION

We demonstrated the applicability of our proposed method by applying the method on two datasets, specifically a dataset of rental properties advertisements and Persian Wikipedia pages related to the city of Tehran, Iran. Divar, a popular platform in Iran that is used by people to advertise their properties to sale, were crawled in a period of one month. The dataset contains 11393 records of advertisements. Each record includes a title and a description of the property in Persian language, pricing information, construction year, the area of the property, the name of the neighborhood where it is located, and locational data. In case of Wikipedia pages, we utilized SPARQL on Wikidata Query Service endpoint to extract articles related to places within the city of Tehran. The dataset contains 623 articles, all of them having locational information and written in Persian language. To clean the dataset, those properties for which the description is less than 10 words were eliminated. This helps in reducing the problem of feature vectors' sparsity. After data cleaning, the dataset contained 5446 records of advertisements. In order to preprocess the descriptions, we employed Stanza [49], a python natural language package, which has been pre-trained for Persian language. First, the documents were segmented into sentences, and then the sentences were tokenized into words. By using POSProcessor module of Stanza, the part of speech of each word was identified and excessive words were removed. To reduce the volume of data for processing, only nouns and adjectives in sentences were considered. In addition, only those words that were pointed out in the whole corpus for

TABLE 1. The accuracy of classification computed for five words in spaces with different number of dimensions and indices, Divar dataset, (the words are translated from Persian to English).

| Index | TF-IDF | | | PPMI | | | |
|-----------------------------|--------------|------|------|------|------|------|------|
| | # dimensions | 5 | 10 | 20 | 5 | 10 | 20 |
| Word 1: <i>luxurious</i> | | 0.42 | 0.49 | 0.52 | 0.44 | 0.52 | 0.62 |
| Word 2: <i>quiet</i> | | 0.26 | 0.35 | 0.41 | 0.26 | 0.36 | 0.40 |
| Word 3: <i>apartment</i> | | 0.31 | 0.34 | 0.36 | 0.37 | 0.42 | 0.46 |
| Word 4: <i>stylish</i> | | 0.39 | 0.44 | 0.48 | 0.45 | 0.51 | 0.66 |
| Word 5: <i>furnished</i> | | 0.35 | 0.39 | 0.46 | 0.42 | 0.49 | 0.55 |

TABLE 2. The accuracy of classification computed for five words in spaces with different number of dimensions and indices, Wikipedia dataset, (the words are translated from Persian to English).

| Index | TF-IDF | | | PPMI | | | |
|--------------------------------|--------------|------|------|------|------|------|------|
| | # dimensions | 5 | 10 | 20 | 5 | 10 | 20 |
| Word 1: <i>historical</i> | | 0.57 | 0.61 | 0.63 | 0.58 | 0.63 | 0.63 |
| Word 2: <i>park</i> | | 0.67 | 0.69 | 0.69 | 0.65 | 0.67 | 0.68 |
| Word 3: <i>neighborhood</i> | | 0.34 | 0.35 | 0.36 | 0.39 | 0.42 | 0.42 |
| Word 4: <i>central</i> | | 0.51 | 0.57 | 0.59 | 0.53 | 0.57 | 0.60 |
| Word 5: <i>building</i> | | 0.46 | 0.51 | 0.53 | 0.46 | 0.49 | 0.50 |

at least 20 times were considered. This ensures the strength of the meaning found in the future steps and enhances the accuracy of classification [50]. Each document was assigned with a BoW containing its nouns and adjectives. In order to construct the feature vectors, both TF-IDF and PPMI indices were calculated and the results of each weighting approach were compared. Conceptual spaces with 5, 10, and 20 dimensions were generated by applying MDS on dissimilarity matrix resulting from angular differences of vectors. A SVM classifier with a linear kernel was trained to locate a hyper-plane on the space. Then, the coefficients of the equation of the hyper-plane, as the vector representing the meaningful direction, was computed. The results of the above procedure for five selected words are presented in Table 1 and Table 2, respectively for Divar dataset and Wikipedia articles. The accuracy assessment in the tables measures the ratio of the number of correctly classified points to the number of all points.

As can be seen from Table 1 and Table 2, documents in higher dimensional spaces are grouped more efficiently such that the classifier can divide the points more accurately. A comparison between the spaces induced from TF-IDF and PPMI implies that by leveraging PPMI, higher accuracies are achieved.

B. COMPARISON WITH BASELINE MODELS

The word embedding resulting from the proposed method was compared with BERT and word2vec as the baseline

TABLE 3. The comparison of the classification task for Divar dataset, performed by the proposed method against baseline models.

| Model | Accuracy | Micro-average F1-score |
|--|----------|------------------------|
| <i>ParsBERT</i> | 72.70 | 62.25 |
| <i>Word2vec</i> | 69.28 | 59.24 |
| <i>Directional Conceptual Spaces (Divar dataset)</i> | 74.25 | 62.48 |

TABLE 4. The comparison of the classification task for Wikipedia dataset, performed by the proposed method against baseline models.

| Model | Accuracy | Micro-average F1-score |
|--|----------|------------------------|
| <i>ParsBERT</i> | 66.36 | 59.52 |
| <i>Word2vec</i> | 65.48 | 57.70 |
| <i>Directional Conceptual Spaces (Wikipedia dataset)</i> | 67.25 | 60.14 |

models widely used in research and practice. For the case of BERT, we used ParsBERT, which is based on BERT architecture and is pre-trained for the Persian language. Table 3 and Table 4 summarize the comparisons made among the models. As shown in the tables, for the case of Divar dataset, the proposed method outperforms word2vec by about 3 percent, in terms of the micro-average F1-score. In terms of F1-score, the results of the proposed method are comparable with that of the ParsBERT model. However, in terms of the accuracy of the classification, the proposed method outperforms both word2vec and BERT models by about 5 percent and 1.5 percent, respectively. In case of Wikipedia dataset, the proposed method yields higher accuracies than ParsBert and word2vec in the scale of 1.3 percent and 2.7 percent, respectively. Again, in terms of F1-score, the proposed method outperforms ParsBert and word2vec by 1.0 percent and 4.2 percent, respectively.

C. EVALUATION

For each word listed in Table 1 and Table 2, an item containing the word was selected and three most similar items were computed. As the proposed approach identifies meaningful directions, a mere Euclidean distance could not guarantee the similarity of two items. Hence, a metric that identifies similarity with respect to the found direction is required. To this aim, the document vectors are projected onto the meaningful direction. Given an item’s vector \mathbf{X} and a meaningful direction \mathbf{S}_w , by vector calculus, the projection of \mathbf{X} on \mathbf{S}_w is computed as Eq. (8).

$$X' = \left(\frac{\mathbf{X} \cdot \mathbf{S}_w}{|\mathbf{S}_w|^2} \right) \mathbf{S}_w \tag{8}$$

By projecting the items on the meaningful directions, the calculation of similarity between two items is reduced to subtracting their corresponding projections (Fig. 3).

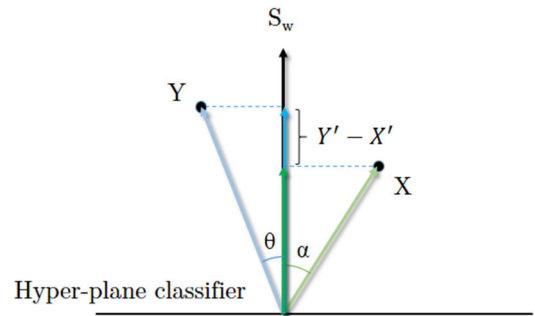


FIGURE 3. Magnetization The proposed method of computing similarity of items in the conceptual spaces.

TABLE 5. Shannon entropy values for five selected topics found in Divar dataset.

| Topic No. | Topic | Shannon entropy |
|-----------|--|-----------------|
| #1 | {luxurious, stylish, Jacuzzi, roof-garden, well-off} | 5.24 |
| #2 | {quiet, tree, alley} | 16.14 |
| #3 | {apartment, parking lot, elevator} | 12.85 |
| #4 | {city center, accessibility, store} | 3.22 |
| #5 | {rebuilt, accessories} | 9.70 |

The more two items are similar, the more the magnitude of the subtraction vector is close to zero. Another advantage of leveraging this approach over using a Euclidean distance is that the similarity achieved through this approach is directional. This characteristic is crucial in recommender systems so that the results can easily be adapted to the users’ needs. Since the users of a place-based recommender system are often seeking items similar to their desired item, it is not interesting to recommend items that are located far from their ideal location. In other words, as we are dealing with place-based recommendations, similarities cannot ensure spatial correlation of items. Therefore, high similarity is not a suitable criterion for recommendation in such scenarios. We employed Shannon entropy to find out how much a word implies spatial correlation. Shannon entropy calculates how much the outcome of a probability distribution would be dispersed. The lower the entropy, the lower the amount of dispersion. Shannon entropy is calculated as Eq. (9).

$$H(p) = - \sum_1^N p_i \log(p_i) \tag{9}$$

We interpret the similarity values of items as probability distributions over words. Table 5 and Table 6 outline the selected items, the normalized similarity values for three most similar items, and entropy values of each word.

Table 5 indicates that the words *quiet*, *tree*, *alley*, *apartment*, *parking*, *lot*, *elevator*, and *apartment* are not good choices for place-based recommendation. This is intuitive as the items associated with these terms do not refer to a specific region and are often more dispersed in the space. The results shown in Table 6 are even more interesting as the topics are much easier interpretable and the terms in a topic are more

TABLE 6. Shannon entropy values for five selected topics found in Wikipedia dataset.

| Topic No. | Topic | Shannon entropy |
|-----------|---------------------------------|-----------------|
| #1 | {historical, central, old} | 3.07 |
| #2 | {park, area, hectare} | 10.16 |
| #3 | {metro, square, street} | 20.12 |
| #4 | {district, neighborhood, place} | 22.70 |
| #5 | {university, school, institute} | 13.45 |

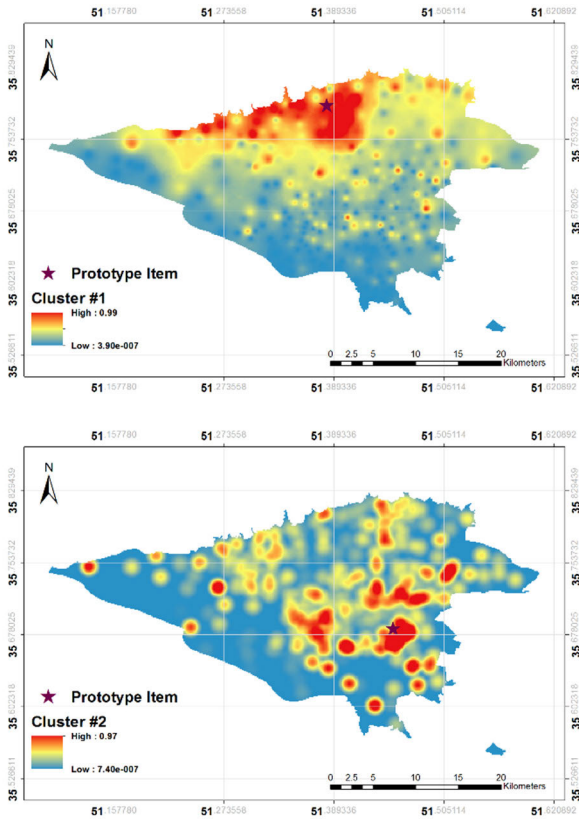


FIGURE 4. Similarity maps of a low-entropy topic (top) and a high-entropy topic (bottom) for Divar dataset.

coherent. For instance, Topic #5 vividly refers to Wikipedia articles about the universities and famous schools in Tehran. A comparison between the two tables shows the importance of using entropy to measure the geo-indicativeness of terms. Although the topics found in Wikipedia articles are semantically more coherent, the entropies are generally high. This shows that the terms of these topics rarely describe their underlying neighborhood or space. Topic #1 in Table 6 is an exception, because they are probably extracted from articles related to historical neighborhoods or buildings in those regions. To evaluate the proposed entropy approach, we sketched the similarity values of items on the map of Tehran. Figures 4 and 5 show the interpolated maps of similarities for Clusters #1 and #2 of Table 5 and Cluster #1 and #4 of Table 6, respectively.

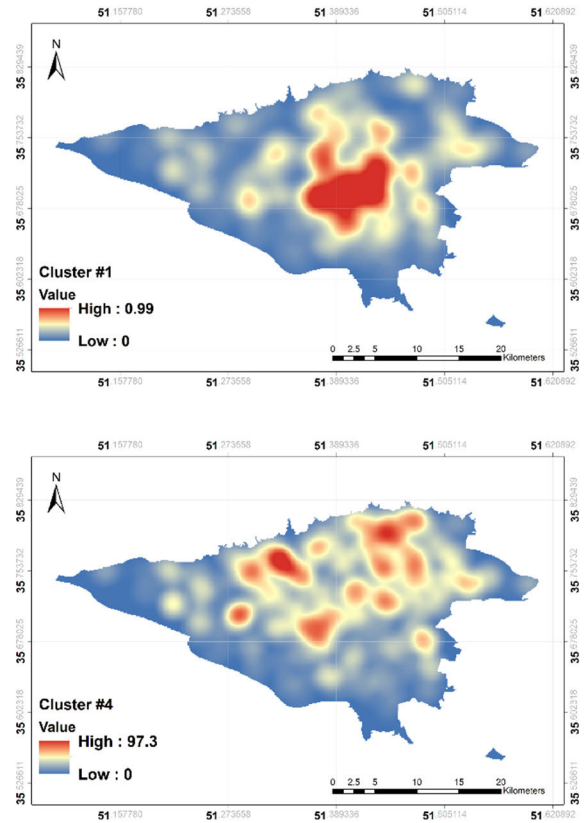


FIGURE 5. Similarity maps of a low-entropy topic (top) and a high-entropy topic (bottom).

In addition, to evaluate the validity of entropies, we compared the results of low-entropy and high-entropy topics of Divar dataset. In order to evaluate the results, we depicted the prototype item of the clusters to 20 individuals (residents of Tehran) and asked them to indicate whether they expected the recommended items or not. They rated the results on a base of 0 to 1, with 0 indicating a completely irrelevant result and 1 indicating a very similar result. Fig. 6 demonstrates the comparison of the average users’ judgments against the computed values.

As can be seen in the plot, the ratings of users to the items of Cluster #1 satisfyingly match the computed similarities. The average error of the computed similarities is 11%. By a correlation coefficient of 0.79, the plot also illustrates that the trend of the similarity assignments matches the users’ ratings. However, for Cluster #2, the average error of results is 23%, which stresses the role of entropy in identifying similar items in place-based recommendations.

The proposed method inherently has the potential to be used in any location-based recommender system with items of textual nature. In this method, there is no presupposition regarding the use case of the recommender system. However, our method collects features through texts associated with place-based items. Considering that in a place recommender system, users expect to find items physically close to their desired item, the quality of the recommendation mainly

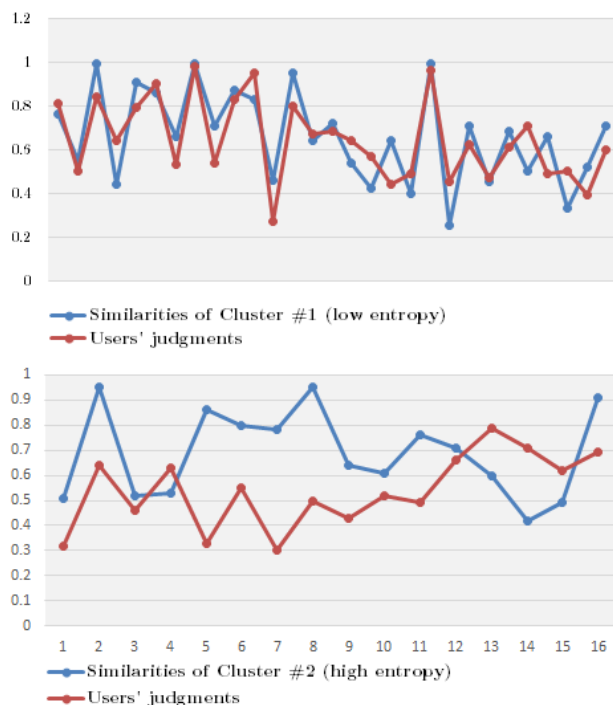


FIGURE 6. The comparison of users' judgments against computed similarities for a high-entropy (top) and a low-entropy (bottom) topic.

depends on the quality of the documents. In other words, the more documents contain features related to specific geographic regions, the less entropy and better recommendations.

VI. CONCLUSION

In this paper, we proposed a method of extracting semantic information based on embeddings and conceptual spaces, which is useful in finding items similarity in recommender systems. The proposed method includes leveraging term saliency indices in the documents and inferring a conceptual space by applying MDS technique on the vector space. Then, a linear classifier is used to divide the conceptual space based on the presence of a given word in the documents. By identifying the perpendicular direction to this hyper-plane, a meaningful direction in the conceptual space is extracted. The projection of items' vectors on the meaningful direction helps in computing the similarities among items while considering the extracted semantic. To apply the methodology, we collected a textual dataset of rental properties in Tehran, Iran, through a web crawling procedure. In addition, to prove the generality of the method, we applied it on a dataset of about 600 Persian Wikipedia articles. The results show that, in some certain number of dimension and up, the classifier can perfectly distinguish between those documents containing a given word and the documents lacking it. The advantage of using the proposed approach is the use of common, yet not complicated, ML techniques. By applying this approach, we transformed the documents to a conceptual space which is interpretable and from which semantic information can be

extracted. Moreover, the approach used here to compute the dissimilarity is directional, which helps in customizing the recommender systems to the users' needs. The evaluation results show that the average accuracy of computing item similarity with the proposed approach is 88%.

REFERENCES

- [1] N. Passalis and A. Tefas, "Learning bag-of-embedded-words representations for textual information retrieval," *Pattern Recognit.*, vol. 81, pp. 254–267, Sep. 2018.
- [2] M. Al-Snadi, M. Al-Ayyoub, Y. Jararweh, and O. Qawasmeh, "Enhancing aspect-based sentiment analysis of Arabic hotels' reviews using morphological, syntactic and semantic features," *Inf. Process. Manage.*, vol. 56, no. 2, pp. 308–319, 2019.
- [3] X. Chen, S. Wang, Y. Tang, and T. Hao, "A bibliometric analysis of event detection in social media," *Online Inf. Rev.*, vol. 43, no. 1, pp. 29–52, Feb. 2019.
- [4] T. Di Noia, C. Magarelli, A. Maurino, M. Palmonari, and A. Rula, "Using ontology-based data summarization to develop semantics-aware recommender systems," in *Proc. Eur. Semantic Web Conf.*, 2018, pp. 128–144.
- [5] S. Qassimi and E. H. Abdelwahed, "The role of collaborative tagging and ontologies in emerging semantic of web resources," *Computing*, vol. 101, no. 10, pp. 1489–1511, Oct. 2019.
- [6] Z. Geng, G. Chen, Y. Han, G. Lu, and F. Li, "Semantic relation extraction using sequential and tree-structured LSTM with attention," *Inf. Sci.*, vol. 509, pp. 183–192, Jan. 2020.
- [7] L. Nizzoli, M. Avvenuti, M. Tesconi, and S. Cresci, "Geo-semantic-parsing: AI-powered geoparsing by traversing semantic knowledge graphs," *Decis. Support Syst.*, vol. 136, pp. 313–346, 2020.
- [8] G. Bakal, P. Talari, E. V. Kakani, and R. Kavuluru, "Exploiting semantic patterns over biomedical knowledge graphs for predicting treatment and causative relations," *J. Biomed. Informat.*, vol. 82, pp. 189–199, Jun. 2018.
- [9] B. Jang, I. Kim, and J. W. Kim, "Word2Vec convolutional neural networks for classification of news articles and tweets," *PLoS ONE*, vol. 14, no. 8, Aug. 2019, Art. no. e0220976.
- [10] D. Richter, S. Winter, K.-F. Richter, and L. Stirling, "Granularity of locations referred to by place descriptions," *Comput., Environ. Urban Syst.*, vol. 41, no. 8, pp. 88–99, 2013.
- [11] D. Richter, S. Winter, K.-F. Richter, and L. Stirling, "How people describe their place: Identifying predominant types of place descriptions," in *Proc. 1st ACM SIGSPATIAL Int. Workshop Crowdsourced Volunteered Geograph. Inf.*, Nov. 2012, pp. 30–37.
- [12] Y.-F. Tuan, "Space and place: humanistic perspective," in *Philosophy in Geography*. Cham, Switzerland: Springer, 1979, pp. 387–427.
- [13] J. Portugali, "Complexity theory as a link between space and place," *Environ. Planning A, Economy Space*, vol. 38, no. 4, pp. 647–664, Apr. 2006.
- [14] B. Helleland, "Place names and identities," *Oslo Stud. Lang.*, vol. 4, no. 2, pp. 1–22, Jul. 2012.
- [15] K. Stock, C. B. Jones, S. Russell, M. Radke, P. Das, and N. Aflaki, "Detecting geospatial location descriptions in natural language text," *Int. J. Geograph. Inf. Sci.*, vol. 36, no. 3, pp. 547–584, Mar. 2022.
- [16] K. Stock, R. C. Pasley, Z. Gardner, P. Brindley, J. Morley, and C. Cialone, "Creating a corpus of geospatial natural language," in *Proc. Int. Conf. Spatial Inf. Theory*, 2013, pp. 279–298.
- [17] B. Li and L. Han, "Distance weighted cosine similarity measure for text classification," in *Proc. Int. Conf. Intell. Data Eng. Automated Learn.*, 2013, pp. 611–618.
- [18] P. D. Turney and P. Pantel, "From frequency to meaning: Vector space models of semantics," *J. Artif. Intell. Res.*, vol. 37, pp. 141–188, Feb. 2010.
- [19] F. Gasparetti, D. Gavalas, S. Ilarri, F. Ricci, and Z. Yu, "Mining social networks for local search and location-based recommender systems," *Pers. Ubiquitous Comput.*, vol. 23, no. 2, pp. 179–180, Apr. 2019.
- [20] B. AlBanna, M. Sakr, S. Moussa, and I. Moawad, "Interest aware location-based recommender system using geo-tagged social media," *ISPRS Int. J. Geo-Inf.*, vol. 5, no. 12, p. 245, Dec. 2016.
- [21] G. Cai, K. Lee, and I. Lee, "Itinerary recommender system with semantic trajectory pattern mining from geo-tagged photos," *Expert Syst. Appl.*, vol. 94, pp. 32–40, Mar. 2018.
- [22] N. J. Yuan, Y. Zheng, L. Zhang, and X. Xie, "T-finder: A recommender system for finding passengers and vacant taxis," *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 10, pp. 2390–2403, Oct. 2013.

- [23] Y. Lai, Z. Lv, K.-C. Li, and M. Liao, "Urban traffic Coulomb's law: A new approach for taxi route recommendation," *IEEE Trans. Intell. Transp. Syst.*, vol. 20, no. 8, pp. 3024–3037, Aug. 2019.
- [24] P. Kefalas and Y. Manolopoulos, "A time-aware spatio-textual recommender system," *Expert Syst. Appl.*, vol. 78, pp. 396–406, Jul. 2017.
- [25] S. Ojagh, M. R. Malek, S. Saeedi, and S. Liang, "A location-based orientation-aware recommender system using IoT smart devices and social networks," *Future Gener. Comput. Syst.*, vol. 108, pp. 97–118, Jul. 2020.
- [26] T. Zuva, "Image content in location-based shopping recommender systems for mobile users," *Adv. Comput., Int. J.*, vol. 3, no. 4, pp. 1–8, Jul. 2012.
- [27] S. Loh, F. Lorenzi, R. Saldaña, and D. Licthnow, "A tourism recommender system based on collaboration and text analysis," *Inf. Technol. Tourism*, vol. 6, no. 3, pp. 157–165, Jan. 2003.
- [28] E. Asani, H. Vahdat-Nejad, and J. Sadri, "Restaurant recommender system based on sentiment analysis," *Mach. Learn. With Appl.*, vol. 6, pp. 100–114, Dec. 2021.
- [29] J. Capdevila, M. Arias, and A. Arratia, "GeoSRS: A hybrid social recommender system for geolocated data," *Inf. Syst.*, vol. 57, pp. 111–128, Apr. 2016.
- [30] S. Amara and R. R. Subramanian, "Collaborating personalized recommender system and content-based recommender system using TextCorpus," in *Proc. 6th Int. Conf. Adv. Comput. Commun. Syst. (ICACCS)*, Mar. 2020, pp. 105–109.
- [31] J. Zhao, J. Wang, M. Sigdel, B. Zhang, P. Hoang, M. Liu, and M. Korayem, "Embedding-based recommender system for job to candidate matching on scale," 2021, *arXiv:2107.00221*.
- [32] X. Tao, N. Sharma, P. Delaney, and A. Hu, "Semantic knowledge discovery for user profiling for location-based recommender systems," *Hum.-Centric Intell. Syst.*, vol. 1, nos. 1–2, p. 32, 2021.
- [33] P. Bafna, D. Pramod, and A. Vaidya, "Precision based recommender system using ontology," in *Proc. Int. Conf. Energy, Commun., Data Anal. Soft Comput. (ICECDS)*, Aug. 2017, pp. 3153–3160.
- [34] M. Riyahi and M. K. Sohrabi, "Providing effective recommendations in discussion groups using a new hybrid recommender system based on implicit ratings and semantic similarity," *Electron. Commerce Res. Appl.*, vol. 40, Mar. 2020, Art. no. 100938.
- [35] L. O. Colombo-Mendoza, R. Valencia-García, A. Rodríguez-González, G. Alor-Hernández, and J. J. Samper-Zapater, "RecomMetz: A context-aware knowledge-based mobile recommender system for movie showtimes," *Expert Syst. Appl.*, vol. 42, no. 3, pp. 1202–1222, Feb. 2015.
- [36] F. Vasile, E. Smirnova, and A. Conneau, "Meta-Prod2Vec: Product embeddings using side-information for recommendation," in *Proc. 10th ACM Conf. Recommender Syst.*, Sep. 2016, pp. 225–232.
- [37] J. Liu, S. Zheng, G. Xu, and M. Lin, "Cross-domain sentiment aware word embeddings for review sentiment analysis," *Int. J. Mach. Learn. Cybern.*, vol. 12, no. 2, pp. 343–354, Feb. 2021.
- [38] M. Faruqui, J. Dodge, S. K. Jauhar, C. Dyer, E. Hovy, and N. A. Smith, "Retrofitting word vectors to semantic lexicons," 2014, *arXiv:1411.4166*.
- [39] P. Gärdenfors, *Conceptual Spaces: The Geometry of Thought*. Cambridge, MA, USA: MIT Press, 2004.
- [40] J. Aisbett and G. Gibbon, "A general formulation of conceptual spaces as a meso level representation," *Artif. Intell.*, vol. 133, nos. 1–2, pp. 189–232, Dec. 2001.
- [41] P. Gärdenfors and F. Zenker, "Editors' introduction: Conceptual spaces at work," in *Applications of Conceptual Spaces*. Cham, Switzerland: Springer, 2015, pp. 3–13.
- [42] M. Raubal, "Formalizing conceptual spaces," in *Proc. 3rd Int. Conf. Formal Ontol. Inf. Syst.*, 2004, pp. 153–164.
- [43] B. Adams and M. Raubal, "A metric conceptual space algebra," in *Proc. Int. Conf. Spatial Inf. Theory*, 2009, pp. 51–68.
- [44] J. Derrac and S. Schockaert, "Inducing semantic relations from conceptual spaces: A data-driven approach to plausible reasoning," *Artif. Intell.*, vol. 228, pp. 66–94, Nov. 2015.
- [45] H. Banaee, E. Schaffernicht, and A. Loutfi, "Data-driven conceptual spaces: Creating semantic representations for linguistic descriptions of numerical data," *J. Artif. Intell. Res.*, vol. 63, pp. 691–742, Nov. 2018.
- [46] J. A. Bullinaria and J. P. Levy, "Extracting semantic representations from word co-occurrence statistics: A computational study," *Behav. Res. Methods*, vol. 39, no. 3, pp. 510–526, Aug. 2007.
- [47] W. C. Gartner, "Tourism image: Attribute measurement of state tourism products using multidimensional scaling techniques," *J. Travel Res.*, vol. 28, no. 2, pp. 16–20, Oct. 1989.
- [48] M. C. Hout, M. H. Papesch, and S. D. Goldinger, "Multidimensional scaling," *Wiley Interdiscipl. Rev., Cogn. Sci.*, vol. 4, pp. 93–103, Jan. 2013.
- [49] P. Qi, Y. Zhang, Y. Zhang, J. Bolton, and C. D. Manning, "Stanza: A Python natural language processing toolkit for many human languages," 2020, *arXiv:2003.07082*.
- [50] M. Karimi, M. S. Mesgari, and R. S. Purves, "A comparative assessment of machine learning methods in extracting place functionality from textual content," *Trans. GIS*, vol. 26, no. 8, pp. 3225–3252, Dec. 2022.



OMID R. ABBASI received the B.Sc. degree in geomatics engineering from the University of Isfahan, Iran, in 2014, and the M.Sc. degree in geospatial information systems from the K. N. Toosi University of Technology, Tehran, Iran, where he is currently pursuing the Ph.D. degree in geospatial information systems. He is currently studying the formalization of places in GIScience by the theory of conceptual spaces. His research interests include social media mining, geocomputation, and natural language processing. As a web developer, he is interested in geospatial web services and open geospatial consortium standards.



ALI A. ALESHEIKH received the B.Sc. degree in surveying engineering from the K. N. Toosi University of Technology, Tehran, Iran, in 1988, the M.Eng. degree in geomatics engineering from the University of New Brunswick, Fredericton, Canada, in 1993, and the Ph.D. degree in geospatial information systems from the University of Calgary, Canada, in 1998. He is currently a Full Professor of geospatial information systems with the K. N. Toosi University of Technology. He is a supervisor of over 62 M.Sc. and 24 Ph.D. students. His research interest includes modeling and managing uncertainties in object-oriented geospatial information systems.

...