



Cognitive Insights into Sentic Spaces Using Principal Paths

Edoardo Ragusa¹ · Paolo Gastaldo¹ · Rodolfo Zunino¹ · Marco Jacopo Ferrarotti² · Walter Rocchia² · Sergio Decherchi³

Received: 14 December 2018 / Accepted: 10 May 2019 / Published online: 16 July 2019
© Springer Science+Business Media, LLC, part of Springer Nature 2019

Abstract

The availability of an effective embedding to represent textual information is important in commonsense reasoning. Assessing the quality of an embedding is challenging. In most approaches, embeddings are built using statistical properties of the data that are not directly interpretable by a human user. Numerical methods can be inconsistent with respect to the target problem from a cognitive view point. This paper addresses the issue by developing a protocol for evaluating the coherence between an embedding space and a given cognitive model. The protocol uses the recently introduced notion of principal path, which can support the exploration of a high-dimensional space. The protocol provides a qualitative measure of concept distributions in a graphical format, which allows the embedding properties to be analyzed. As a consequence, the tool mitigates the black-box effect that is typical of automatic inference processes. The experimental section involves the characterization of AffectiveSpace, demonstrating that the proposed approach can be used to describe embeddings. The reference cognitive model is the hourglass model of emotions.

Keywords Topological analysis · Sentiment analysis · Affective computing · Concept embedding

Introduction

Concept-level sentiment analysis is a challenging area of research on natural language processing (NLP) [1]. It involves the semantic investigation of text, mostly using inference techniques, to aggregate conceptual and affective information associated with natural language opinions. Concept-level sentiment analysis has countless potential applications in interdisciplinary areas including stock market prediction, political forecasting, social network analysis, social stream mining, and human-robot interaction.

Sentiment analysis is by definition a function of the way that the human brain perceives emotions. As a consequence, in the past few years, different researchers have tried to embed commonsense knowledge in deep learning models. In some cases, sentiment information has been included in embeddings based on deep learning strategies [2] or directly in the models [3–5]. Furthermore, cognitive analysis can target cultural peculiarities related to specific languages

[6], suggesting that cognitive analysis is necessary for development in the field. The proposed methodology provides a graphical inspection tool to simplify the integration of commonsense knowledge.

The quality of the investigative process relies on a suitable conversion of the textual information into a numerical format, generally called embedding. Thus, developing an effective embedding is critical in all branches of text mining [3, 4, 7–11]. An embedding can be formalized as a transformation $F : T \rightarrow R^d$ where $d \in N^+$ and T is the set of admissible fragments of text. There are two macro categories of embedding, namely, general-purpose and task-specific. The first group comprises remapping methodologies based on mere statistical properties of the text. The second group comprises methodologies for building numerical spaces based on task-specific constraints. Task-specific remapping strategies can be obtained by exploiting general-purpose embeddings. In this regard, astonishing results have been achieved by using deep learning to tailor these general-purpose representations to specific tasks using external resources [12–14]. Recently, Li et al. [2] proposed a technique for developing a more refined sentiment lexicon based on the GloVe embedding [15] evolved with the use of external resources. In general, it is computationally demanding and data

✉ Edoardo Ragusa
edoardo.ragusa@edu.unige.it

demanding to derive a task-specific embedding from a general-purpose embedding. For this reason, the use of task-specific mappings is preferred in most real-world scenarios [10, 15]. Various sentiment lexicons, such as the multi-perspective question answering (MPQA) corpus [16], NRC [17], and SenticNet [18] provide sentiment polarities of words directly. Although focused, this scalar information is limited when it comes to modeling the complex interactions between terms and sentiments. Another example of sentiment-oriented embedding can be found in [19], where the authors proposed using the Harvard psychological dictionary and Loughran-McDonald financial sentiment dictionary to build a sentiment space. Concept-based approaches have been more effective than word-based approaches because sentiments are usually related to combinations of many terms. The most prominent implementation of this approach is AffectiveSpace [20], a 100-dimensional embedding of concepts.

Assessing the quality of an embedding has several challenging issues. First, a proper tool is needed to inspect the topology and properties of data distributions in high-dimensional spaces [21]. This issue is typically addressed with dimensionality reduction techniques [22–26], which can also support visual inspection [27]. A second major issue is to define a proper metric to measure the similarity of concepts in the embedding space. Standard methodologies use clustering algorithms that adopt either the Euclidean distance or a normalized scalar product. The performance of these approaches is strongly dependent on correctly defining a similarity measure that can capture the intrinsic structure of the data [28] and on fruitfully defining the cluster [29]. Finally, one must check the consistency between the concept distribution in the embedding space and the psychological models. Accordingly, one must demonstrate that concepts that are similar in the embedding space are indeed similar in terms of human perception.

Although measuring the embedding quality is difficult and ambiguous, it is necessary for the development of the field. A qualitative inspection of the embedding allows designers to select the most promising approaches in the early stages and to validate the proposed solution in the advanced stages. Moreover, it reduces the so-called black box problem that affects most data-driven strategies.

This paper contributes to the definition of a novel protocol for assessing the quality of an embedding. The principal path concept [30] is used to explore and characterize the high-dimensional space that arises from an embedding. A principal path is a path that connects two points in data space and that tries to pass through the *local mid of the data*. A principal path is a local version of the principal curve algorithm [31], in which the end points are additionally constrained to specific sample points. These paths have a cognitive value because, by passing through the

center of the data (in a local way), they try to capture the most abstract morphing path between the two user-selected end points. To compute principal paths, [30] proposed a regularized version of the popular k -means (hereafter referred as RKM) algorithm. We note that the link between k -means and RKM is mainly formal, in that the aim of the two tools is completely different. The ability to provide information about data topology in high-dimensional spaces is particularly appealing when it comes to analyzing the quality of an embedding. This is because it can also reveal details of the transitions from one concept to another. This information cannot be obtained with a clustering algorithm.

In the proposed protocol, the ability to find principal paths with the RKM algorithm in a high-dimensional affective space is critical to evaluating the consistency between the distribution of concepts in this space and a reference psychological model. Thus, given a set of concepts, a reference model for characterizing the affective information conveyed by each concept, and the concept embedding, the goal is to directly and visually analyze the concept distribution topology, which is expected to be coherent with the reference psychological model. Here, we have developed such an analytical instrument, which allows the qualitative study of cognitive reasoning in high-dimensional spaces derived from concept embedding. The proposed method has been tested on AffectiveSpace, graphically representing the concept distribution in this embedding and confirming the consistency between AffectiveSpace and Plutchik's theories [32]. The resulting graphical representation is based on the peculiarities of the RKM algorithm, which explores the transition between concepts. As a consequence, the space is not explored with clustering. This is important because, by testing the transitions between the true reference concepts, we can consistently validate the psychological model from a cognitive point of view. In addition, the paper presents a qualitative characterization of AffectiveSpace, where the user is provided with a series of cognitive descriptors about the data distribution.

The rest of the paper is organized as follows. The “**Preliminaries**” section briefly reviews the presented approach. The “**Proposed Methodology**” section formalizes the proposed protocol for assessing the coherence between an embedding and a reference cognitive model. The “**Experimental Results**” section presents and analyzes the results of the experimental campaign. Finally, some concluding remarks are given in the “**Conclusion**” section.

Preliminaries

This section briefly reports on the tools and concepts underlying the work. The “**Hourglass Model**” section analyzes the hourglass model, which is the psychological model

selected to represent the sentimental content of concepts. The “AffectiveSpace” section introduces AffectiveSpace, i.e., the embedding that will be analyzed. Finally, the “Principal Paths in Data Space by Regularized Kernel *k*-Means” section introduces the path-finding method at the core of the proposed inspection protocol.

Hourglass Model

The Hourglass of Emotions [33] is a suitable tool for characterizing concepts in terms of human emotions. It is based on Plutchik’s studies on human emotions and involves reorganizing sentiments around four independent dimensions whose different levels of activation make up the total emotional state of the mind. The Hourglass of Emotions is depicted in Fig. 1. The basic idea is that different independent resources determine the emotional state. Different conditions arise from turning on one set of these resources and turning off another set. In each configuration, the behavior of the brain changes. For

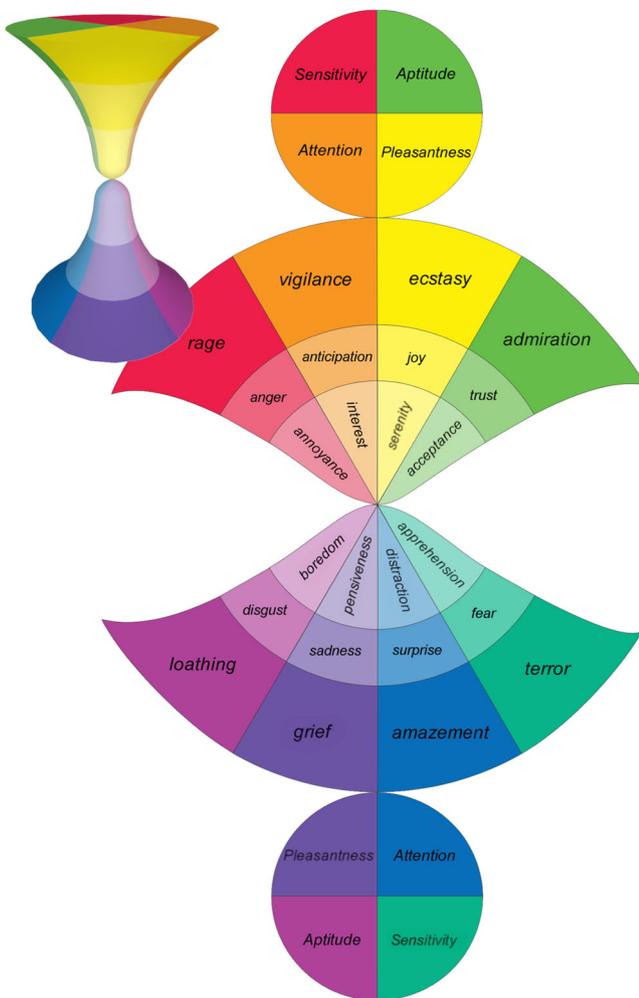


Fig. 1 The Hourglass of Emotions [20]

example, the state of “anger” appears to privilege a set of resources that support an immediate reaction to external stimuli, while also suppressing a set of resources that make us act more prudently. The hourglass shape reflects the idea that an emotion can be identified only if it is strong enough, i.e., we cannot feel a specific emotion like “fear” or “amazement” without that emotion being reasonably strong. The model does not merely classify the affective states into basic emotional categories. Rather, the affective state is described in four concomitant, yet independent dimensions: *sensitivity*, *aptitude*, *attention*, and *pleasantness*. Each dimension is characterized by six levels of activation, which determine the intensity of the expressed/perceived emotion as an analog value ranging from $[-1, +1]$. For example, as depicted in Fig. 1, the six activation levels of *sensitivity* are rage, anger, annoyance, apprehension, fear, and terror. The six activation levels for each of the affective dimensions provide a labelled set of 24 basic emotions, allowing the model to specify the affective information associated with text in a dimensional form and a discrete form.

In this work, the SenticNet tool [18] provides the knowledge base that associates a set of natural language concepts with their numerical description in the four-dimensional space derived from the Hourglass of Emotions.

AffectiveSpace

AffectiveSpace [20] is a compressed version of AffectNet, which is an affective commonsense knowledge base developed using the graphic representation of the Open Mind corpus (called ConceptNet [34]) and WordNet-Affect [35] (a linguistic resource for lexically representing affect). AffectNet is a semantic network where multi-word expressions of commonsense knowledge are nodes, and the links between these are relations between concepts (see Fig. 2 for an example).

A matrix representation of AffectNet is obtained by dividing each assertion into two parts: the first part is

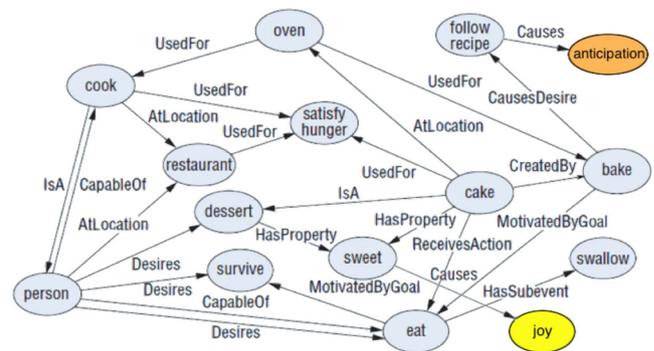


Fig. 2 Example of AffectNet structure for the concept *cake* [36]

the concept, and the second part is the assertion with the first or second concept left unspecified (e.g., “is a kind of liquid” or “a wheel is part of”). Numerical values are obtained by considering the reliability of the assertions. Positive or negative numbers are associated with single assertions based on reliability, with a magnitude that grows logarithmically with the confidence score. Finally, to visualize the concept-relation-concept structure of the graph in a matrix format, the data are reshaped, with every known concept of some statement being a row, and every known semantic feature (relationship + concept) being a column (as per Fig. 3). One advantage of matrix representation is the possibility of performing cumulative analogy. This is executed by selecting a set of nearest neighbors (in terms of similarity) of the input concept, and then projecting known properties of this set onto unknown properties of the concept.

The matrix representation of AffectNet involves thousands of columns. This creates computational issues that limit its use for many practical applications. To address this, Cambria et al. proposed AffectiveSpace [20]. AffectiveSpace is an embedding for concepts, which is built using random projections of the matrix representation of AffectNet. Its purpose is to compress the semantic features associated with commonsense concepts and thus to better perform analogical reasoning on these. AffectiveSpace collects 50,000 concepts, which are remapped in a 100-dimensional space. This high-dimensional space is at the core of the Sentic Computing Engine [37]. Although used in many applications, there is little topological information about this model. Furthermore, all the available descriptions are based on the pair {angle, module} of each concept only.

Principal Paths in Data Space by Regularized Kernel *k*-Means

Principal paths are paths that connect data in a feature space and that are designed to pass locally through the mid of the

data [30]. Principal paths are a local version of principal curves [31], where the end points are constrained to fixed positions. To compute a principal path, a regularized version of *k*-means was introduced in [30] (RKM). RKM is an algorithm for selecting smooth paths in data space. It is based on the idea that, given a dataset and two reference points in a vector space $X \in R^d$ where $d \in N^+$, a morphing between them can be found by exploiting the available data. This morphing is described as a discrete path, comprising a set of waypoints that are selected based on the data manifold. Although RKM is, in functional terms, a regularized version of the *k*-means, it is not a clustering algorithm. Its purpose is different. The solution of *k*-means coincides with the set of centers that better describe the data population. In contrast, the solution of RKM is the morphism between two points. This conceptual difference means that the two methodologies have different applications.

Formally speaking, consider a set of points $X = \{x_j \in R^d\}$, $j = 1, \dots, N$ and two points w_0 and $w_{N_c+1} \in R^d$. The path connecting the two points w_0 and w_{N_c+1} is described as an ordered set W of N_c waypoints (prototypes) $w \in R^d$.

Algorithmically, RKM can be seen as a function of two extreme points w_0 and w_{N_c+1} , a set of points in the same space X , and a discretization parameter N_c , with an output consisting of an ordered set of points W :

$$W = \text{RKM}(w_0, w_{N_c+1}, X, N_c); \tag{1}$$

The path is found by minimizing the standard *k*-means cost function with the addition of a regularization term, which considers the distance between pairs of ordered waypoints. The cost function can be formalized as:

$$\min_W \frac{\gamma}{2} \sum_{i=1}^N \sum_{j=1}^{N_c} \|x_i - w_j\|^2 \delta(u_i, j) + \frac{\lambda}{2} \sum_{i=0}^{N_c} \|w_{i+1} - w_i\|^2 \tag{2}$$

where $\delta(u_i, j)$ gives the membership of the *i*th sample to the *j*th cluster/waypoint. This novel cost function comprises

Concepts	Semantic Features (relationship+concept)							
	..	Causes joy	IsA event	UsedFor housekeeping	LocatedAt party_venue	PartOf celebration	MotivatedByGoal clean_room	..
⋮		⋮	⋮	⋮	⋮	⋮	⋮	
wedding	..	0.94	0.86	0	0.79	0.88	0	..
broom	..	0	0	0.83	0	0	0.87	..
buy_cake	..	?	0.78	0	0.80	0.91	0	..
birthday	..	0.97	0.85	0	0.99	0.98	0	..
sweep_floor	..	0	0	0.79	0	0	0.91	..
⋮		⋮	⋮	⋮	⋮	⋮	⋮	

Fig. 3 Example of AffectNetMatrix structure [36], where rows represent concepts, and columns contain the numerical values of the features

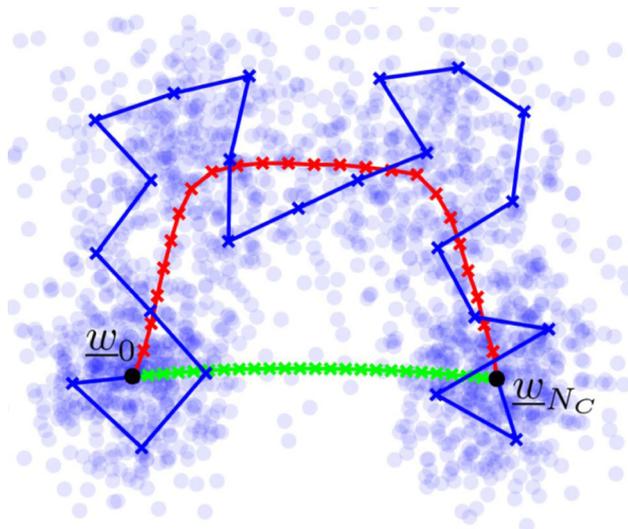


Fig. 4 Example of hyper-parameter influence in the shape of the path

two terms weighted by the hyper-parameters γ and λ (they can be reduced to one, but we follow the original formulation where evidence maximization was used):

$$\Omega(\mathbf{W}, \mathbf{u}, \mathbf{X}, \gamma, \lambda) = \gamma \Omega_X(\mathbf{W}, \mathbf{u}, \mathbf{X}) + \lambda \Omega_W(\mathbf{W}) \quad (3)$$

The first term coincides with the standard k -means cost function. The second term induces a path topology due to the centroid ordering and controls the level of smoothness of the path. Figure 4 provides a graphical example of the algorithm behavior in a two-dimensional space, for different values of the regularization hyper-parameters. In the plot, data are represented as blue dots, and waypoints as crosses. The blue line refers to a configuration where the first cost function term is prominent. The green line refers to a configuration where the second term is prominent. The red line refers to a configuration with the right trade-off between the two, and that corresponds to a principal path.

The minimum of the cost function is obtained with an expectation maximization algorithm [38]. The procedure can be summarized as:

- *E-step*: consider iteration t of the optimization process and \mathbf{W}_t the set of prototypes. Minimizing $\Omega(\mathbf{W}, \mathbf{u}, \mathbf{X}, \gamma, \lambda)$ with respect to \mathbf{u} is equivalent to minimizing $\Omega_X(\mathbf{W}, \mathbf{u}, \mathbf{X})$ that is the standard k -means cost function.

$$u_{i,t+1} \leftarrow \underset{j}{\operatorname{arg\,min}} \| \mathbf{x}_i - \mathbf{w}_{j,t} \|^2 \quad (4)$$

- *M-step*: based on the newly computed \mathbf{u}_{t+1} , minimization of $\Omega(\mathbf{W}, \mathbf{u}, \mathbf{X}, \gamma, \lambda)$ with respect to \mathbf{W} is performed.

The closed-form solution of the M-step can be formalized by defining the following quantities:

- The $N_C \times d'$ boundary condition matrix

$$\mathbf{B}_{i,\cdot} = \begin{cases} \mathbf{w}_0 & \text{if } i = 1 \\ \mathbf{w}_{N_C+1} & \text{if } i = N_C \\ 0 & \text{otherwise} \end{cases} .$$
- The $N_C \times d'$ centroid matrix: $\mathbf{C}_{i,\cdot} = \sum_{j=1}^N \mathbf{x}_j \delta(u_i, j)$.
- The $N_C \times N_C$ Toeplitz matrix representing the hessian of the regularization term in the cost function:

$$\mathbf{A}_W = \begin{bmatrix} 1 & -\frac{1}{2} & 0 & \dots & 0 \\ -\frac{1}{2} & 1 & -\frac{1}{2} & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & -\frac{1}{2} & 1 & -\frac{1}{2} \\ 0 & \dots & 0 & -\frac{1}{2} & 1 \end{bmatrix} \quad (5)$$

- The $N_C \times N_C$ hessian matrix of the standard k -means cost function, $\mathbf{A}_{X_{i,j}} = \begin{cases} |w_i| & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$

The complete RKM procedure is provided in the pseudocode 1.

This new cost function has been paired with the framework of Bayesian evidence maximization to select the hyper-parameters. Interestingly, the choice of the regularization parameter induced by the evidence framework has the further advantage of virtually removing the other degree of freedom, namely the number of waypoints. It finds almost the same manifold even when varying their number, reducing that parameter to a mere discretization parameter of the path.

Algorithm 1 RKM pseudocode.

RKM Pseudocode:

Input

- \mathbf{X} data matrix
- \mathbf{W}_0 initial waypoints matrix s regularization parameter

Output

- \mathbf{W}_t centroids matrix
- \mathbf{u}_t labels

Procedure RKM($\mathbf{X}, \mathbf{W}_0, s$)

```

 $u_{i,0} = \operatorname{arg\,min}_j \| \mathbf{x}_i - \mathbf{w}_{j,0} \|^2$ 
 $t = 0$ 
do
     $\mathbf{W}_{t+1} = (\gamma \mathbf{A}_X(\mathbf{u}_t) + \lambda \mathbf{A}_W)^{-1} (\gamma \mathbf{C} + \frac{\lambda}{2} \mathbf{B});$ 
     $u_{i,t+1} = \operatorname{arg\,min}_j \| \mathbf{x}_i - \mathbf{w}_{j,t+1} \|^2;$ 
     $t = t + 1;$ 
while  $u_{i,t} \neq u_{i,t-1}$ 
return  $\mathbf{W}_t, \mathbf{u}_t$ 
    
```

Proposed Methodology

This paper aims to provide a protocol to investigate, from a cognitive point of view, the structure of high-dimensional spaces that are induced by an embedding. The methodology must (1) inspect the geometrical configuration of data (i.e., concepts) in the embedding space and (2) link this information to cognitive descriptors. The final goal is to facilitate a visual analysis of the investigation outcomes. Therefore, information should be provided in a way that is intuitive, effective, and concise.

The protocol relies on a set C of n concepts, which are expected to provide a suitable sampling of the high-dimensional space that is induced by the embedding F under analysis. Every concept included in this set should come with a cognitive description. In the present work, this description is based on the hourglass model (as per “Hourglass Model”). As a result, each concept $c_i \in C$ comprises a four-dimensional numerical description h_i . Indeed, let $\{x_1, \dots, x_n\} \in R^d$ be the projections of the n concepts in the embedding space.

The actual design of the protocol can be divided into two steps. First, the RKM algorithm is used to inspect the topology of the concept distribution in the embedding space. Second, the topology is linked to the cognitive descriptors, thus providing an analytical instrument for evaluating the consistency between the geometrical configuration of concepts and the reference cognitive model. The “Space Exploration” section provides details about the implementation of the first step, while the “Descriptors” section addresses the second step.

Space Exploration

The RKM algorithm is a suitable tool for exploring the distribution of the concepts $c_i \in C$ in the embedding via their projections $\{x_1, \dots, x_n\}$. Let $l_{start}, l_{end} \in C$ be two generic concepts that will correspond to the two extremes of the path under analysis. Accordingly, RKM can be used to identify the path that connects l_{start} with l_{end} in R^d . Thus, the algorithm output is the list of intermediate concepts, which characterize the transition induced by the data distribution.

For clarity, it is useful to consider Fig. 4. Here, each point corresponds to the 2D projection x_i of a concept $c_i \in C$. w_0 and w_{N_c+1} are the projections of l_{start} and l_{end} , respectively. The path selected by the algorithm represents the most *natural* transformation between the two concepts. Each waypoint identifies an intermediate step in this transformation. Notably, the N_c prototypes $w_j \in R^d$ placed by RKM are not, in general, coincident with known concepts, i.e., concepts in C . In principle, however, each prototype represents the projection of a known or unknown concept in the embedding space.

Algorithm 2 (Path Generator) formalizes the adopted procedure. Given an embedding F , two concepts belonging to C are used as l_{start} and l_{end} , respectively. Eventually, the two points in R^d corresponding to the concepts l_{start} and l_{end} become w_0 and w_{N_c+1} . The RKM algorithm is used to find the path between the two points, provided by the concept distribution of the embedding. The set cL of N_c prototypes in R^d is then provided as output. The generic prototype w_j is not expected to correspond to the projection of a known concept x_i . Therefore, each single prototype is linked to the closest known concept in terms of Euclidean distance in the remapped space. This metric is coherent with the metric used by the RKM algorithm.

Algorithm 2 Generation of morphism between two known concepts.

Path_Generator;

Input

Embedding F with domain C and relative projections $X \in R^{N \times d}$,
a pair of landmark concepts l_{start} and l_{end} ,
number of prototypes N_c

Output

List of concepts representing transitions cL

Procedure

```

 $w_0 = F(l_{start})$ 
 $w_{N_c+1} = F(l_{end})$ 
 $path = RKM(w_0, w_{N_c+1}, X)$ 
for  $j \leftarrow 1$  to  $length(path)$  do
   $L(j) = \arg \min_{x \in X} \|path(j) - x\|^2$ 
   $cL(j) = F(Lj)^{-1}$ 
end
return  $cL$ ;
```

Algorithm 1 can be applied to any embedding F . Moreover, it can be used to suitably model the topology of a concept distribution in the high-dimensional space induced by an embedding. In this regard, it is convenient to design an exploration strategy that, as landmarks, uses concepts that correspond to the activation of the four dimensions of the hourglass model, i.e., concepts that can be used to obtain a discrete quantization of the space. This setup involves the reasonable assumption that the resulting 24 concepts are included in the embedding under analysis.

The strategy adopted in this work is based on independent analysis of the four dimensions. The goal is to characterize the transitions between the adjacent activations of a given dimension. Thus, let \mathcal{L} be the ordered set of six concepts that characterize a dimension. For example, with *sensitivity* as the dimension to be analyzed, one has $\mathcal{L} = \{\text{rage}, \text{anger}, \text{annoyance}, \text{apprehension}, \text{fear}, \text{terror}\}$. Then, Algorithm 2 is run by using as pairs $\{l_{start}, l_{end}\}$,

respectively, {'rage', 'anger'}, {'anger', 'annoyance'}, {'annoyance', 'apprehension'}, {'apprehension', 'fear'}, and {'fear', 'terror'}. The rationale is that the four dimensions are expected to be independent of each other. Thus, in principle, when moving along one of these dimensions, the activation of the remaining three dimensions should not vary. By linking the five paths that characterize a single dimension, one obtains the transition between the concept that corresponds to the highest activation level and the concept that corresponds to the lowest activation level.

Algorithm 3 (Embedding Exploration) outlines the overall strategy, which leads to 20 paths (five paths for each dimension). As a result, one explores the space by studying the morphisms through the manifolds of each dimension, using as reference points the concepts relative to the activations of the hourglass model. The outcome is not only the information about the morphing between two extremes of a path; rather, the proposed strategy also reveals a set of descriptors (i.e., concepts) for the underlying data manifolds.

Algorithm 3 Complete mechanism of space exploration.

Embedding exploration;

Input

Embedding F with domain C and relative projections $X \in R^{N \times d}$,
 number of prototypes N_c ,
 number of subsampling Q

Output

Found_Paths

Procedure

```

for  $j \leftarrow 1$  to 4 do
   $\mathcal{L} = \{c_{1j}, c_{2j}, c_{3j}, c_{4j}, c_{5j}, c_{6j}\}$ 
  for  $i \leftarrow 1$  to  $\text{length}(\mathcal{L})-1$  do
    for  $q \leftarrow 1$  to  $Q$  do
       $\tilde{C}_q \in C$ 
       $\text{Found\_Paths}_j(i, q) = \text{Path\_Generation}(C$ 
         $= \tilde{C}_q, l_{\text{start}} = c_{i,j}, l_{\text{end}} = c_{i+1,j},$ 
         $N_c = N_c)$ 
    end
  end
end
end
return  $\text{Found\_Paths}_j, j = 1, \dots, 4;$ 

```

Algorithm 3 applies a perturbation scheme to obtain a more complete set of manifold descriptors. Given a pair of concepts $l_{\text{start}}, l_{\text{end}}$, Algorithm 2 is applied Q times. During each run, the set of concepts C is subsampled to obtain a new set of subspaces $\tilde{C}_i \subset C$, with $i = 1, \dots, Q$. At the end of the procedure, one has collected Q different paths, which describe the transitions induced by different

subsamplings of the space. A similar mechanism could have been implemented by using different initializations for the set of prototypes. However, the subsampling procedure can induce more variance in the set of solutions.

Figure 5 summarizes the entire space exploration procedure. Given a dimension, the hourglass model sets the five pairs of concepts to be used as landmarks in the path generation. The given example refers to *sensitivity*. The RKM algorithm is then used to extract a path for each pair of concepts on the embedding space under analysis. The process is repeated Q times, eventually obtaining Q paths (i.e., Q lists of N_c concepts) for each pair of landmarks.

Descriptors

The paths obtained in the first phase, space exploration, do not by themselves provide an analytical instrument for evaluating the consistency of the embedding F with respect to a cognitive model. Each path corresponds to the list of concepts encountered by moving from l_{start} to l_{end} in the embedding space. Indeed, the main goal is to assess the coherence of the cognitive descriptors, which are associated with the concepts that share a common path. Nonetheless, a second goal is to facilitate a direct visual evaluation of the outcomes of this analysis.

This paper proposes two different approaches for achieving these goals. The two approaches differ in the criterion adopted to define the cognitive content of a concept. For each concept encountered along a path, the first approach assigns a tag selected from a predefined group of eight tags. Because the tags are based on the hourglass model (two tags per dimension), one can identify the most active components of the emotional space. The second approach, conversely, assigns four analogical values to each concept. That is, one value is assigned for each principal dimension of the hourglass model. The sections “[Tag descriptors](#)” and “[Distribution descriptors](#)” present the two approaches, respectively, in detail.

Tag descriptors

The first approach tries to model the mechanisms that the human brain applies when perceiving emotion. In principle, it is not possible to identify a specific feeling if it is not reasonably strong. Hence, a concept should be associated with the most prominent component of the feeling under observation.

To this end, the approach relies on eight tags to encode the hourglass model. First, let us quantize each dimension of the hourglass space into two levels: positive values and negative values. Accordingly, one can associate a “positive” tag and a “negative” tag with each dimension. For each dimension, Table 1 gives the tags associated with the two

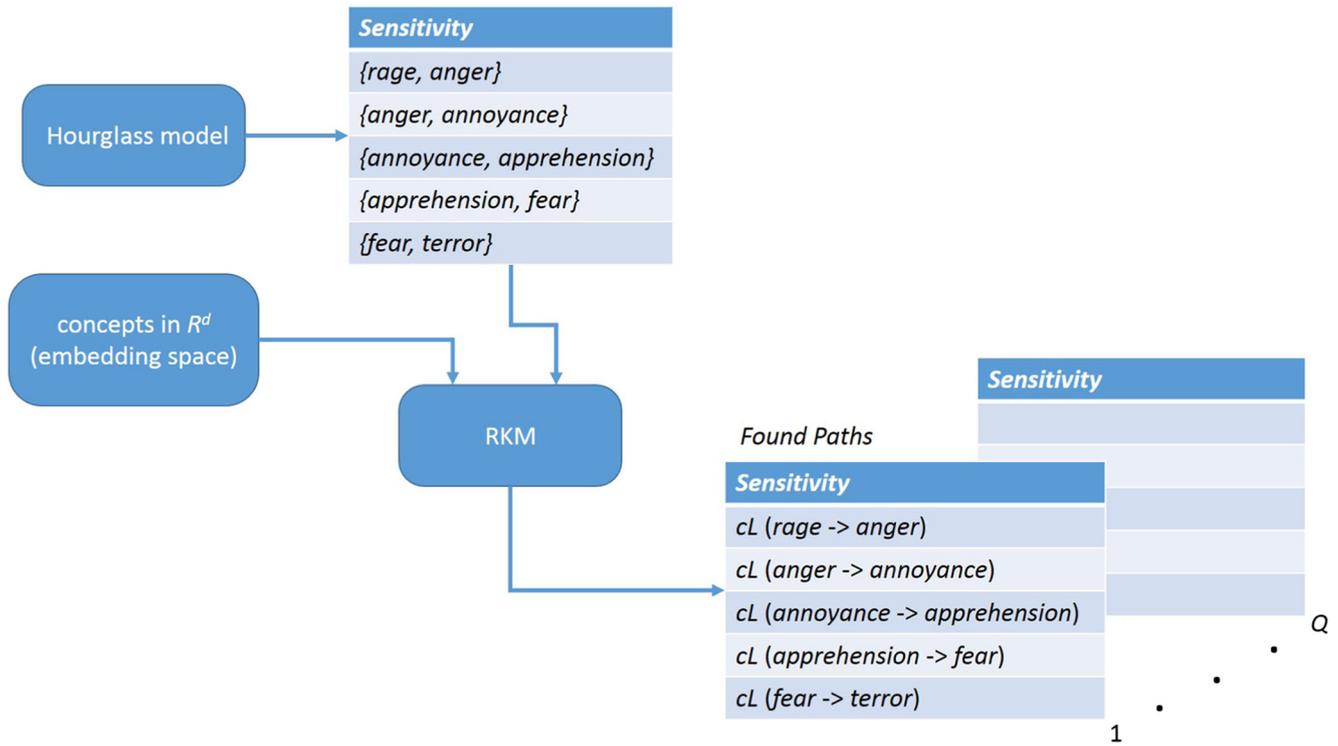


Fig. 5 A schema of the space exploration strategy

levels. Any concept, then, should be associated with one of the eight tags. To achieve this, the following procedure is applied.

Given a concept c_i and the corresponding four-dimensional vector h_i , which provides the activation levels of the hourglass space, one selects the dimension characterized by the highest activation level in absolute value. Thus, should be $h_j = \arg \max_j (|h_i|)$ is the most prominent component. Then, one assigns to c_i the positive tag or the negative tag according to the sign of h_j . For example, if a concept has *aptitude* = 0.3, *sensitivity* = -0.7, *attention* = 0.1, and *pleasantness* = -0.4, then the leading dimension is *sensitivity* and the selected tag is “fear”.

The availability of the pairs {concept,tag} allows the objective characterization of a given path, which includes the N_c concepts encountered while moving from l_{start} to l_{end} in the embedding space. Algorithm 3 provides Q paths for each pair $\{l_{start}, l_{end}\}$. Therefore, the total list of $N_c \times Q$ tags is used. As a result, a path is described in terms of an

eight-dimensional vector t , where each dimension reports the occurrence (expressed in percentage) of a given tag; thus,

$$t_{tag} = \frac{\#_{tag}}{N_c \times Q} \tag{6}$$

where $\#_{tag}$ is the number of times that a specific tag appears in the path.

The descriptor (6) does not take into account the fact that, in principle, the list of concepts included in C is not equally distributed over the eight tags. This affects the expected occurrence of a tag in a given path, thus introducing a bias that might affect the analysis of the properties of a path. Therefore, one may rely on a different descriptor. In this case, the prior probability p_{tag} of a tag is also considered. p_{tag} is the ratio between the number of concepts in C that are associated with the tag and the total number of concepts in C . The prior probability is then used to normalize the descriptor (6). Accordingly:

$$t_{tag} = \frac{\alpha \times \#_{tag}}{N_c \times Q \times p_{tag}} \tag{7}$$

where α is a normalization factor, such that the sum of the percentage of all the tags is 1.

Figure 6 reports the outcome. For each dimension of the hourglass model, the input comprises Q instances of each path, as per Algorithm 3 and Fig. 5. The SenticNet

Table 1 Hourglass model quantization

Name	Positive tag	Negative tag
<i>Aptitude</i>	Admiration	Disgust
<i>Sensitivity</i>	Anger	Fear
<i>Attention</i>	Interest	Surprise
<i>Pleasantness</i>	Joy	Sadness

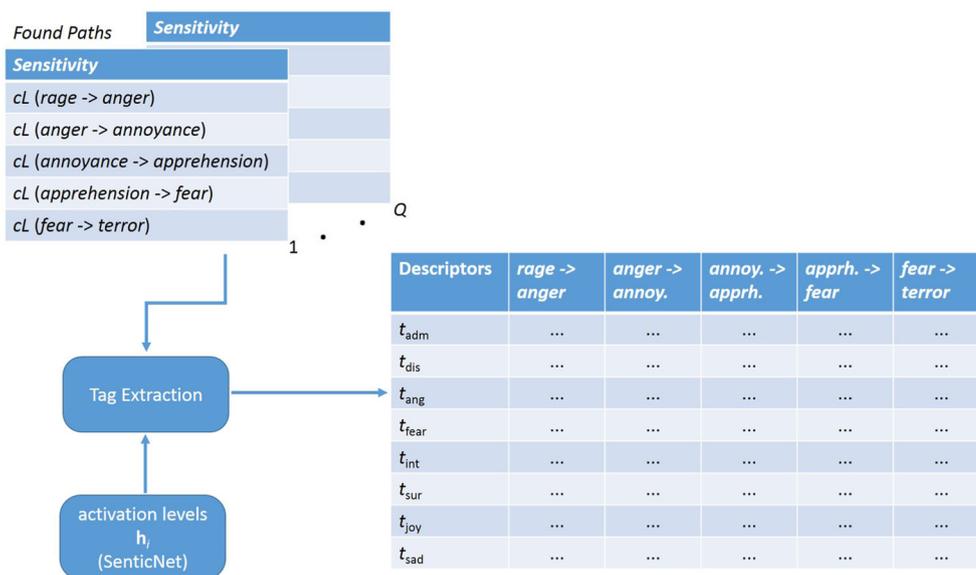


Fig. 6 Block diagram of processing flow for the extraction of tag descriptors

engine provides for each $c_i \in C$ the corresponding vector h_i . For each of the five paths associated with a dimension, one obtains the eight descriptors, computed as per (6) or per (7). Figure 6 reports only the visual representation of the manifold containing the concept related to dimension *sensitivity*. A complete characterization of the space also includes the other three dimensions of the hourglass model.

Distribution descriptors

The second approach exploits a set of descriptors, which provide analogical information about the activation of the hourglass model components. According to Plutchik’s

theory, the four-dimensional representation h_i describes exhaustively the sentiment of a concept c_i . Hence, one can use the activation values to illustrate all the affective information embedded in a single point of the embedding, i.e., a concept.

As stated above, for each dimension, the procedure receives as input the paths generated by Algorithm 3. Given a pair of landmarks $\{l_{start}, l_{end}\}$, the corresponding set of descriptors is generated by processing all the $N_c \times Q$ concepts involved. In this case, the goal is to compute the probability distribution of the six activation levels of each dimension. Thus, a transition from l_{start} to l_{end} is described in terms of four probability distributions, where

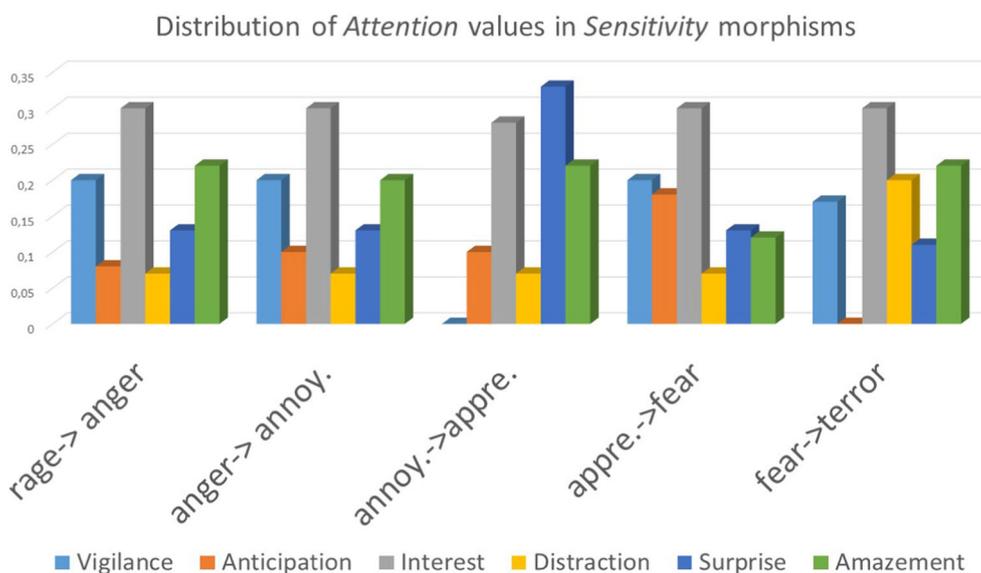


Fig. 7 Distribution of activation values for the attention dimension when analyzing the morphisms inside the sensitivity manifold

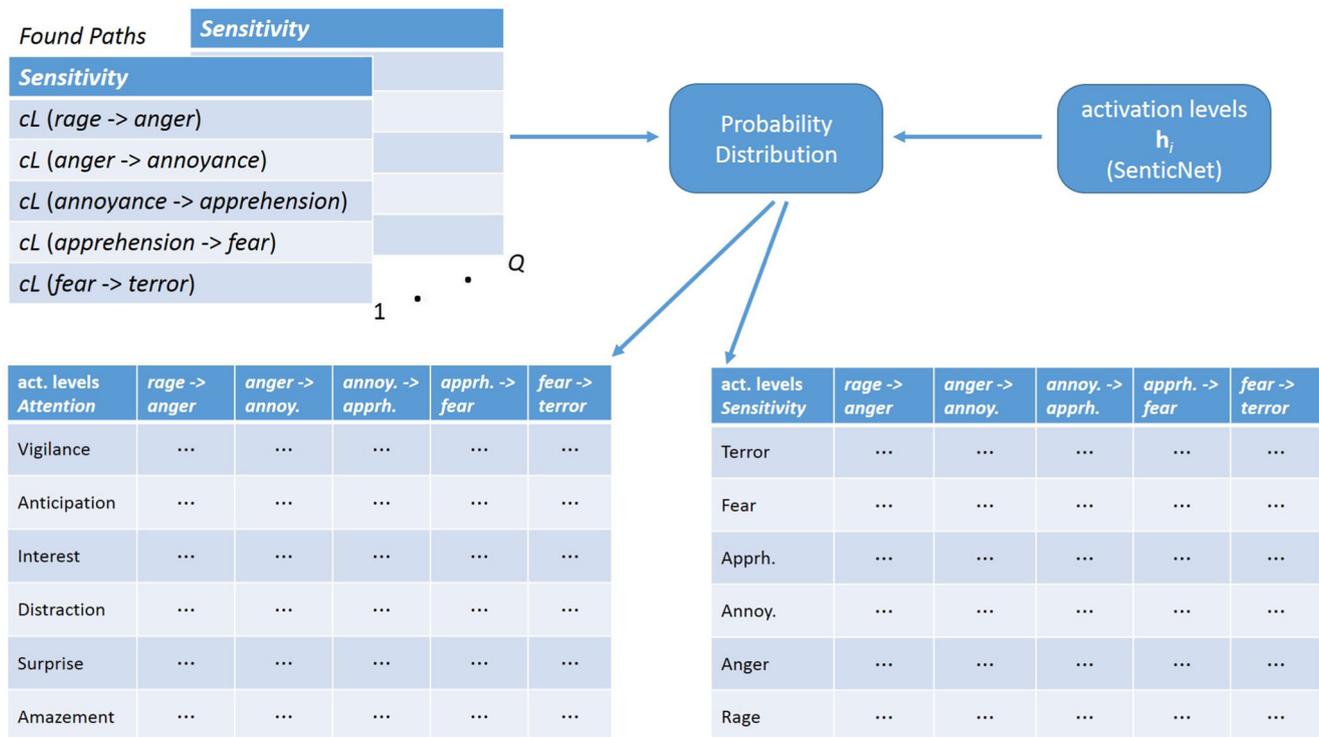


Fig. 8 Block diagram of the processing flow for extracting the distribution descriptors

each distribution involves six bins (one per activation level). Each probability distribution is estimated by exploiting the $N_c \times Q$ concepts involved and their vectors h_i .

The example in Fig. 7 considers the whole morphism between rage and terror, i.e., the morphism between the most positive and the most negative concept in the *sensitivity* dimension. Accordingly, the plot gives the five corresponding transitions on the *x*-axis. For each transition, the plot provides the probability distribution relative to the six activation levels of the *attention* dimension.

Finally, Fig. 8 schematizes the processing flow that leads to the descriptors based on probability distributions. Again, given a dimension of the hourglass model, the processing starts from the Q instances of each path, as per Algorithm 3 and Fig. 5. The SenticNet engine provides for each $c_i \in C$ the corresponding vector h_i . For each pair of landmarks $\{l_{\text{start}}, l_{\text{end}}\}$, the procedure builds the probability distribution of the activation levels for each dimension of the hourglass model. In Fig. 8, the input refers to the *sensitivity* dimension, while the outputs refer, respectively, to probability distributions computed over the activation levels of the *attention* dimension (on the left) and the *sensitivity* dimension (on the right). The procedure generates the analogous descriptors for the probability distributions computed over the activation levels of the *aptitude* and *pleasantness* dimensions, respectively.

Comparison with Existing Methods

The proposed approach directly links vector spaces and cognitive descriptors. To the best of our knowledge, this is the first method of its kind in this regard. Assuming both the descriptors and the psychological model as invariant, the selection of the manifold analysis algorithm is critical. The alternatives include dimensionality reduction techniques, such as T-sne [26] and PCA [23], and manifold learning techniques, such as local linear embedding [24], multi-dimensional scaling [25], and Laplacian eigenmaps [39]. In general, these algorithms try to preserve the global topological structure of the dataset, i.e., they are designed to preserve the local pattern-to-pattern similarity. They thus address a particularly hard problem, especially when dealing with high-distance values. The main difference to principal paths is that principal paths do not try to solve the whole manifold embedding problem. Rather, it directly addresses the problem of locally walking over the manifold, thus solving a different and simpler problem. Being more local, this approach increases the probability of returning a reliable representation. Moreover, clustering techniques such as standard *k*-means do not order the clusters; thus, there is no topological information. They are excellent approaches for defining groups of data, but not for investigating the structure of a manifold.

Table 2 Manifold methods comparison

Method	Aim
Principal Path	Path Finding
T-sne [26]	Projection
PCA [23]	Projection
Local Linear Embedding [24]	Projection
Multi-dimensional scaling [25]	Projection
Laplacian Eigenmaps [39]	Projection
K-means [38]	Clustering

Table 2 summarizes the aforementioned observations.

Experimental Results

This section applies the protocol formalized in the “Proposed Methodology” section to the AffectiveSpace embedding. The goal is to show that the numerical descriptors proposed in the “Descriptors” section are an analytical instrument for investigating the coherence between an embedding space and a reference cognitive model (here, the hourglass model). The experimental characterization of AffectiveSpace is divided into three phases. The first phase focuses on the outcomes of the description based on tags. The second phase analyzes the outcomes of the description based on probability distributions. The third phase addresses the issues around polarization of the results due to the projection of the embedding points in concepts. Overall, the experimental design aims to:

- provide a detailed characterization of AffectiveSpace;
- empirically demonstrate the effectiveness of the proposed methodology in visualizing a cognitive description of the space.

The regularization parameters were obtained using Bayesian evidence maximization and the number of clusters N_c was set to 13. The paths were obtained using the Euclidean distance as metric. Algorithm 3 was run on six subsamplings of the embedding space. Each subsampling included 50,000 concepts (half of the maximum number of available concepts).

Preliminary Analysis

AffectiveSpace comprises 100,000 concepts in a 100-dimensional space. By construction, this configuration is sparse. Figure 9 highlights the issues related to this aspect. This two-dimensional example illustrates a somehow pathological configuration. The concepts are organized in two well-separated clusters. Indeed, the two

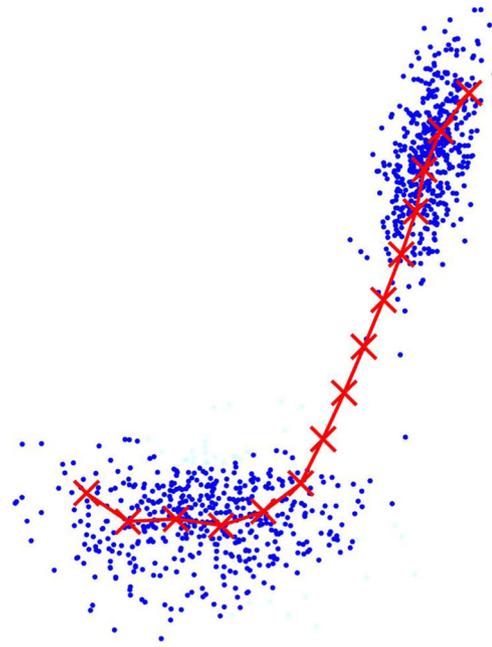


Fig. 9 Example of prototypes in empty regions of the data space

landmarks are placed at the two extreme sides of the distribution. As a result, most of the prototypes (red crosses) describing the path are outside the data distribution. This phenomenon becomes more prominent with increases in the dimensionality. Indeed, it might produce a slight deterioration in the quality of the solution, which is obtained by projecting the prototypes onto the closest concept.

Notably, SenticNet contains the same concepts that are mapped in AffectiveSpace. This solves the problem of possible mismatches between concepts associated with the prototypes and available descriptors. In general, though, a matching strategy should be developed.

The SenticNet engine suffers from an uneven distribution of the tags defined in the “Descriptors” section. Table 3 reports the actual distribution of these tags over the concepts included in SenticNet, which provided the set C of concepts for Algorithm 3. In the table, each row gives a tag and the percentage of concepts associated with that tag with

Table 3 Percentage of tags for the whole dataset

Name	Percentage
admiration	0.0155
disgust	0.0050
anger	0.1237
fear	0.0537
interest	0.1647
surprise	0.0787
joy	0.2893
sadness	0.2696

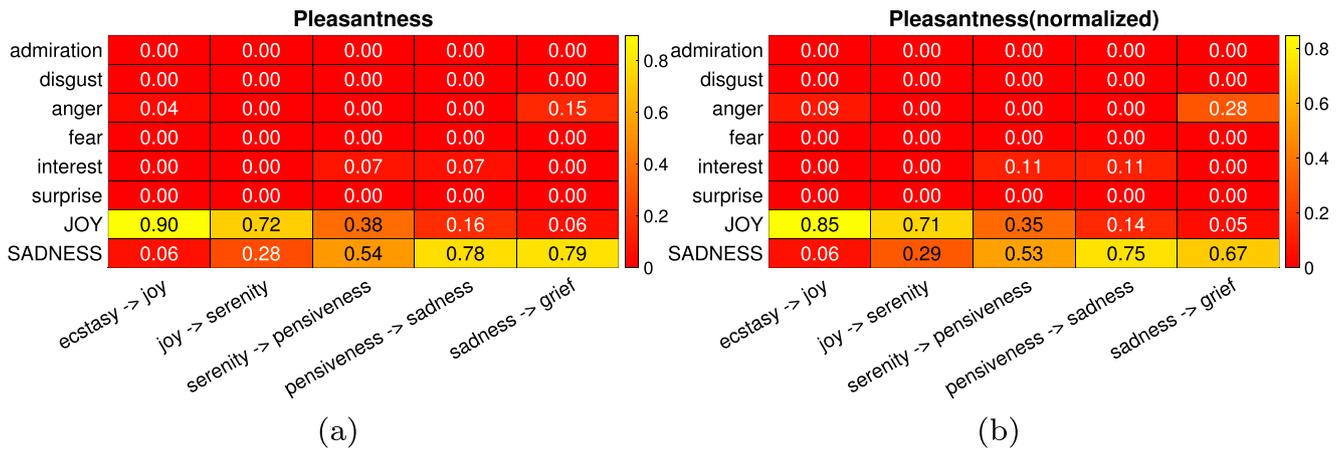


Fig. 10 a, b Pleasantness visualization using tag descriptors

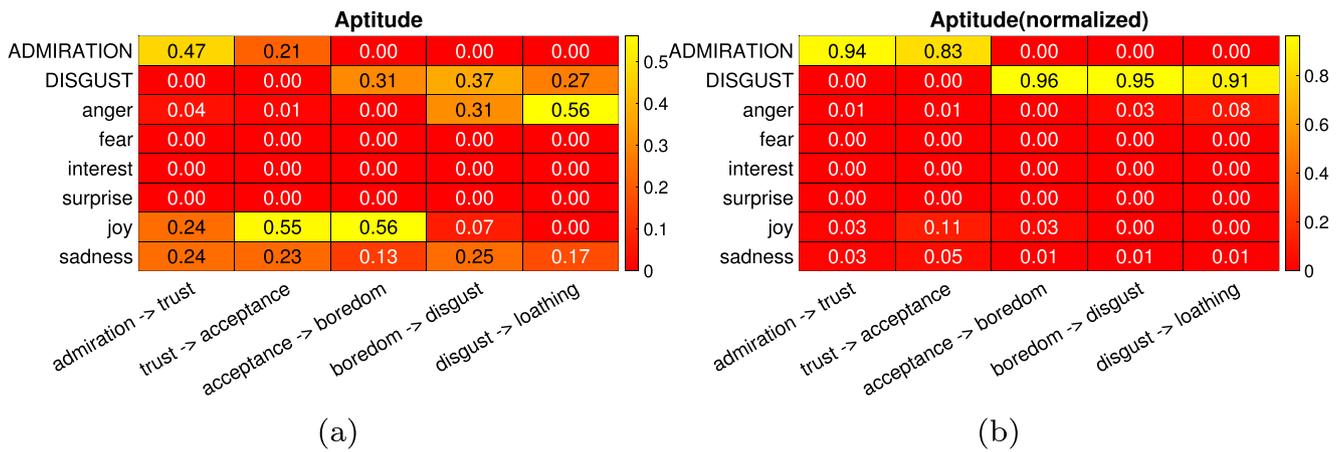


Fig. 11 a, b Aptitude visualization using tag descriptors

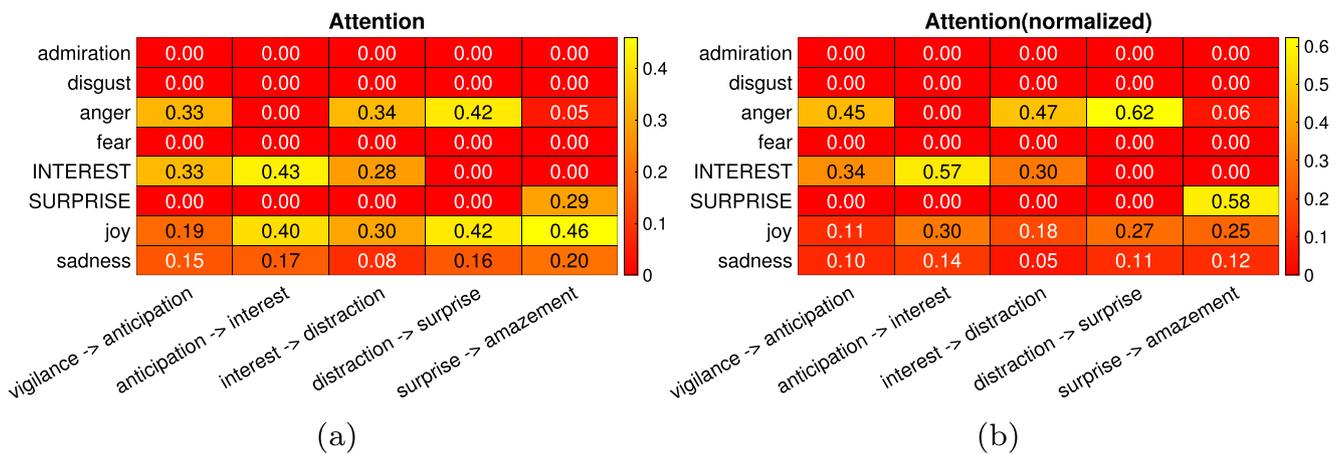


Fig. 12 a, b Attention visualization using tag descriptors

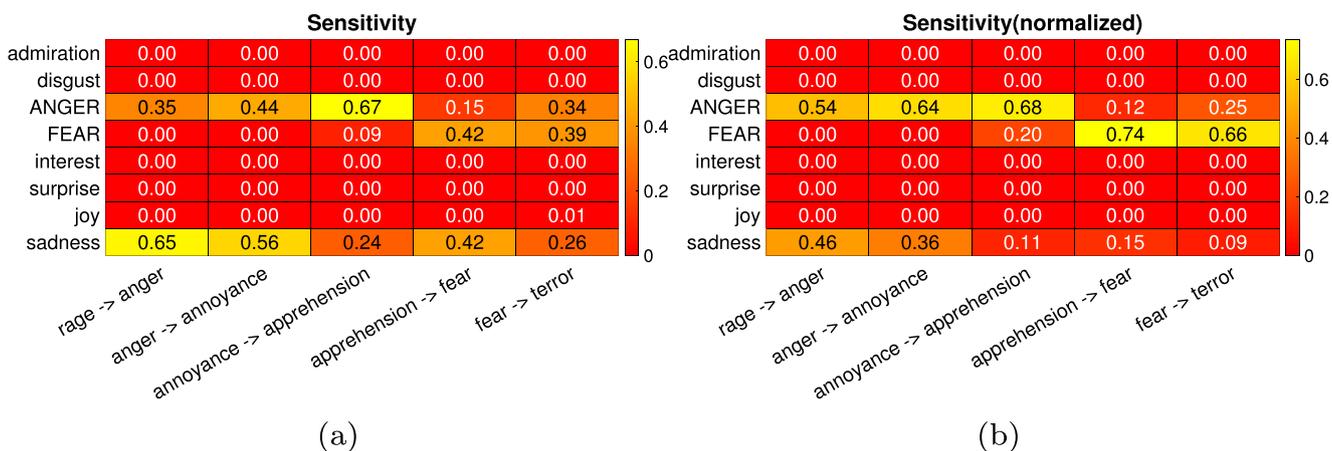


Fig. 13 a, b Sensitivity visualization using tag descriptors

respect to the total number of concepts. Table 3 shows that the probability distribution over the eight tags is far from uniform. Tags joy and sadness comprise almost 50% of the data, while tags admiration and disgust comprise around 2%. The descriptor (6) was designed to address this issue.

Tag Analysis

The first phase of the proposed analysis involves the descriptors introduced in the “Tag descriptors” section. The space has been characterized by exploring the four dimensions of the hourglass model independently, as

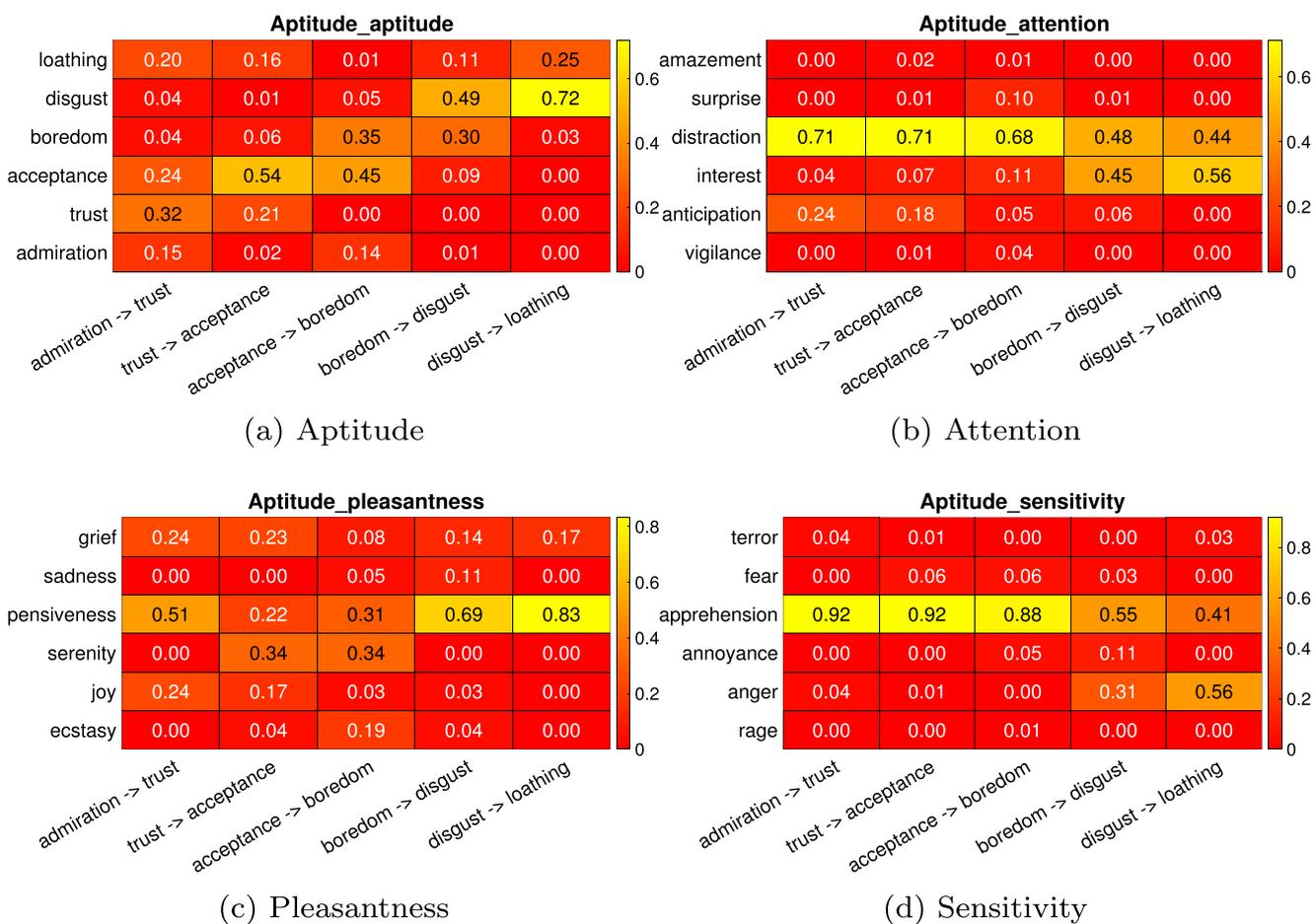


Fig. 14 a–d Aptitude visualization using distribution descriptors

explained in “Space Exploration” section. Figures 10, 11, 12, and 13 provide a visual report of the outcomes obtained by analyzing the paths generated for the *pleasantness*, *aptitude*, *attention*, and *sensitivity* dimensions, respectively. Each figure includes two tables that share the same format. The columns identify the five transitions that characterize the morphism under analysis. Notably, by moving from column 1 to column 5, one inspects that dimension by passing through the corresponding activation levels in the order suggested by the hourglass model. In the table on the left, for each column (i.e., for each transition), the eight rows report on the value assigned to the descriptor (6). Hence, this gives the occurrence (expressed in percentage) of the eight tags defined in Table 1 in the single transition. In the table on the right, the eight rows report on the value assigned to the descriptor (7). The latter table thus introduces a normalization factor to take into account the uneven distribution of tags in the concept list.

In each table, capital letters are used to indicate the names of the two tags that are directly linked to the dimension under analysis (as per Table 1). The goal is to simplify

the interpretation of the plots. In the ideal configuration, we would expect that, for a given dimension, its path in the embedding space would mostly encounter concepts that are linked to the tags expressed by that dimension. This would suggest complete coherence between the cognitive representation of the hourglass model and AffectiveSpace. To further facilitate a prompt visual analysis of the tables, colors are associated with the percentages, which report the tag occurrences. In each table, the corresponding color bar is presented. The color scale runs from red (0%) to yellow, which sets the maximum percentage in the table. Thus, the color bar may have a different numeric range from table to table.

A first outcome of this analysis arises from the tables to the left of each figure (i.e., the tables (a)). When using the non-normalized version of the descriptor, the two tags that represent *pleasantness*, joy, and sadness are active in all four dimensions. That is, a considerable percentage of concepts are associated with those tags, independent of the dimension under analysis. This is a direct consequence of the fact that, in *C*, almost 50% of the concepts are

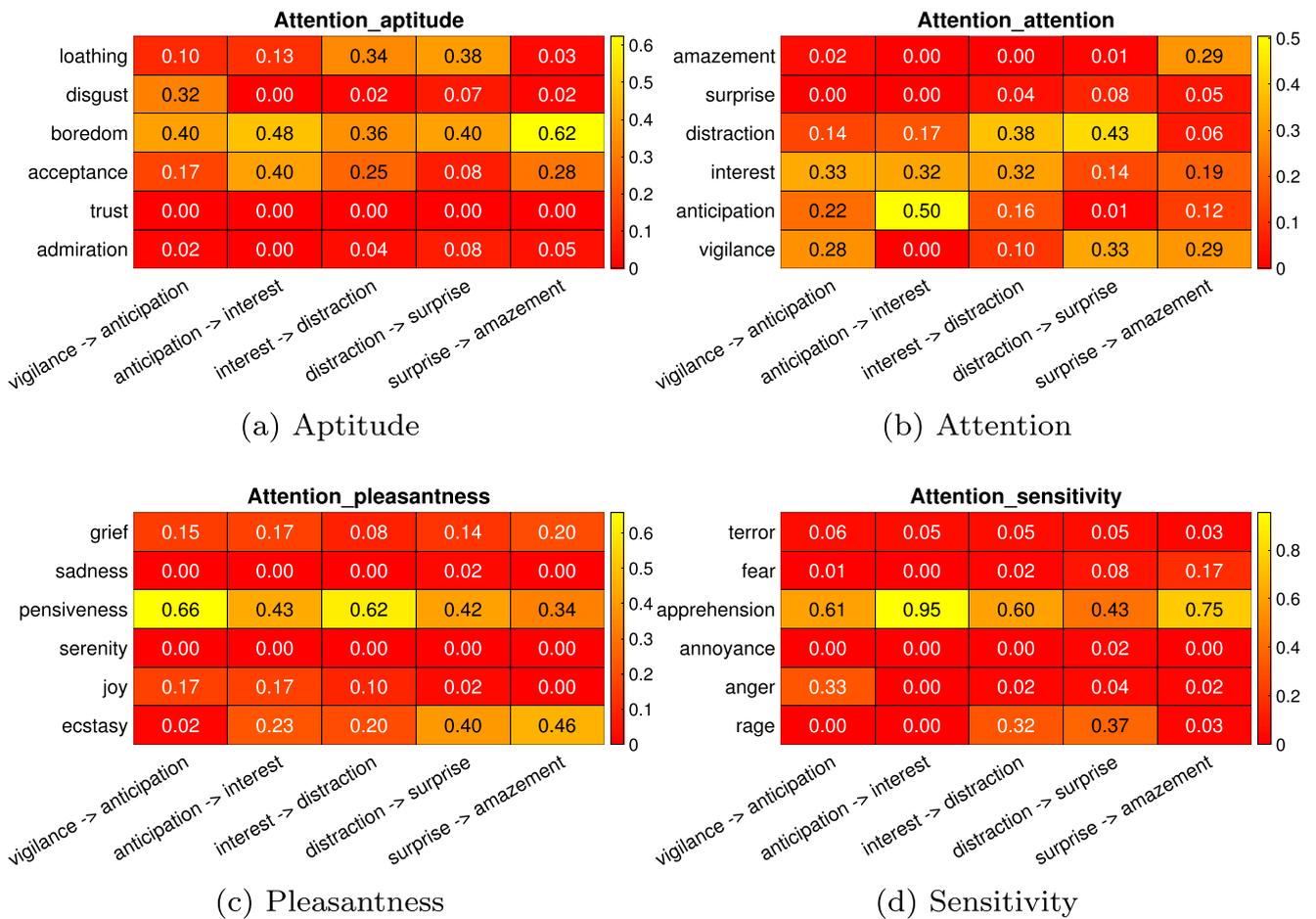


Fig. 15 a–d Attention visualization using distribution descriptors

associated with these tags. Indeed, the tag anger also seems to occur often. Nonetheless, it is possible to see that, for a given dimension, the two corresponding tags (given in capital letters) are always considerably active. Indeed, the positive tag occurs more frequently when the transitions refer to the positive activations (first two columns), while the negative tag occurs more frequently when the transitions refer to the negative activations (last two columns). In fact, this behavior is less evident only in the case of *attention*, where the tags anger, joy, and sadness are also quite active. The tables (b) confirm that, by exploiting the normalization factor, one can obtain a more reliable outcome. In the case of *pleasantness*, *aptitude*, and *sensitivity*, the reference tags are actually the most active. This means that the AffectiveSpace is coherent with the hourglass model. Table (b) in Fig. 12, however, shows that the path relative to *attention* encounters a relevant number of concepts tagged as Anger.

Finally, we briefly comment on the single tables separately. In *pleasantness* (Fig. 10), the tables show that the path mostly encountered concepts tagged as joy or sadness,

as one would expect. Moreover, there is a clear transition from joy to sadness when moving from positive activations to negative activations. Some noise appears in the last path (i.e., from sadness to grief), where 16% of the concepts are associated with anger. From a cognitive point of view, this overlap is reasonable. In the *aptitude* dimension (Fig. 11), table (a) shows the effect of the unbalanced nature of the concept list. The admiration and disgust tags are associated with only 2% of the concepts. Table (b), though, reveals that the list of concepts encountered during the transition is coherent with the cognitive characterization of the dimension. As anticipated above, the *attention* dimension (Fig. 12) seems the least effective in cognitive terms. On the one hand, the interest tag appears as expected when the path goes through the positive activations. On the other hand, the surprise tag occurs only in the last section of the path. Besides, even in table (b), tags related to joy, sadness, and anger are quite active. Finally, the *sensitivity* dimension (Fig. 13) seems to confirm the coherence between the AffectiveSpace and the cognitive representation that is set by the hourglass model.

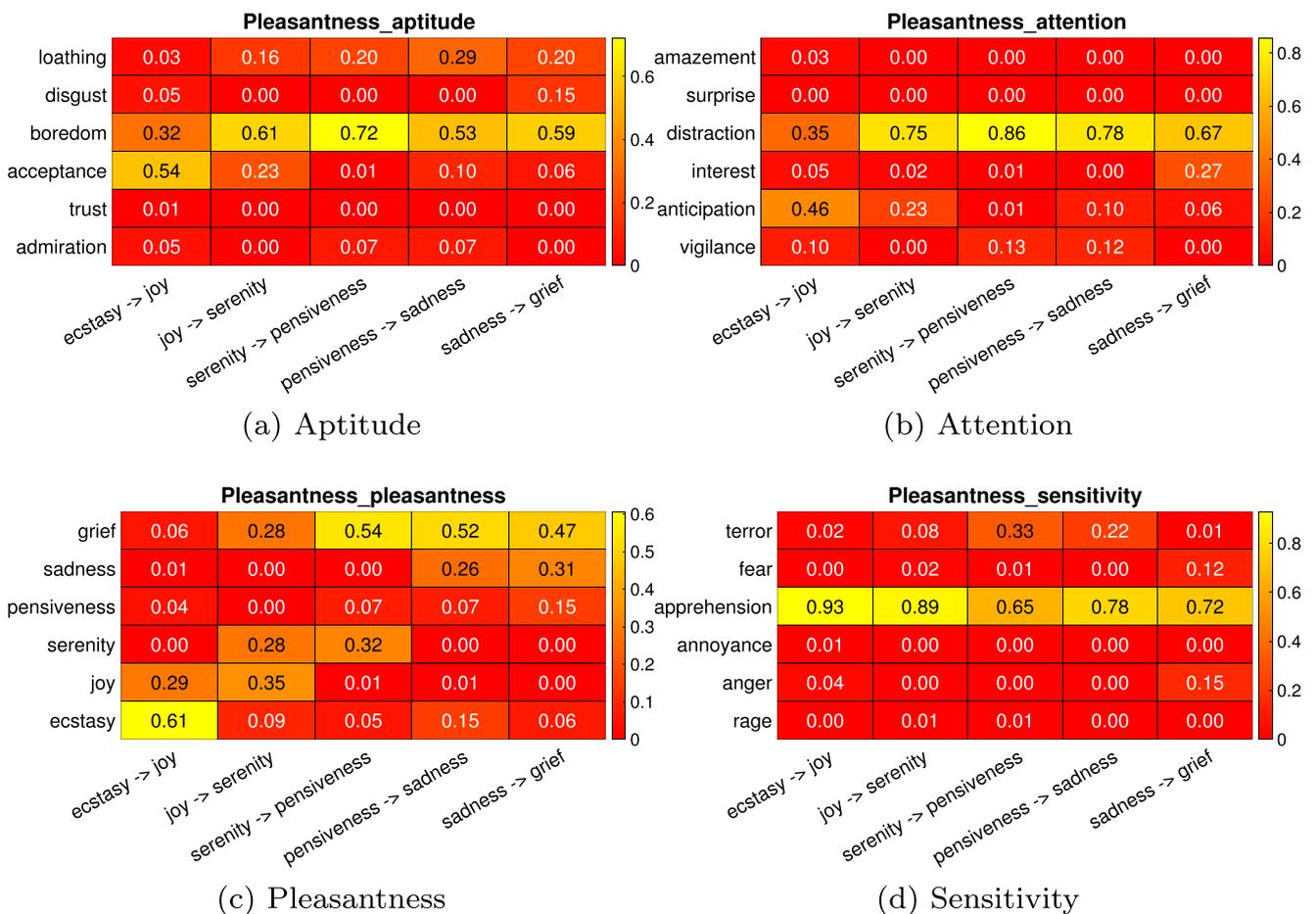


Fig. 16 a–d *Pleasantness* visualization using distribution descriptors

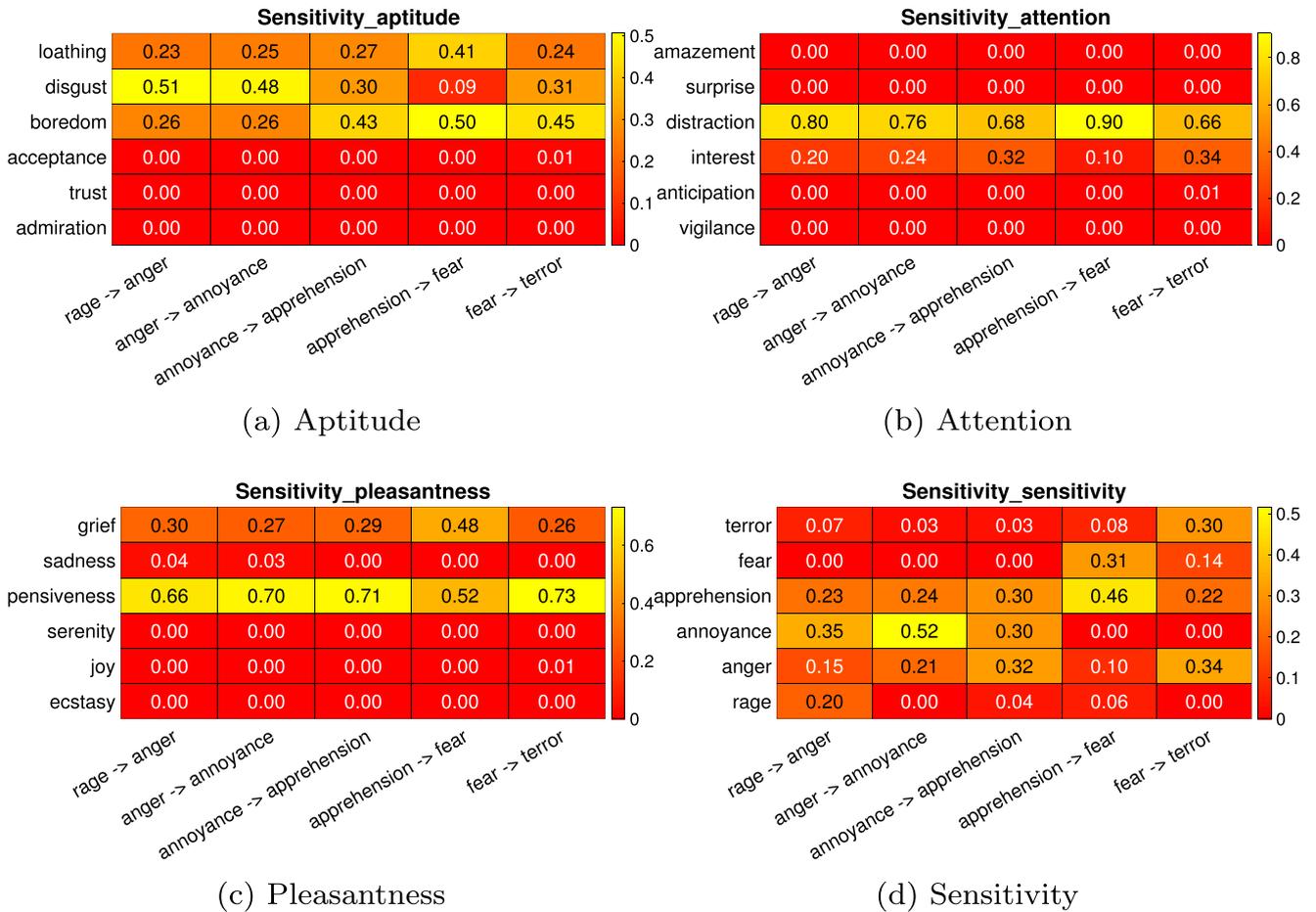


Fig. 17 a–d Sensitivity visualization using distribution descriptors

Distribution Analysis

The second proposed analysis involves the descriptors introduced in the “Distribution descriptors” section. For each dimension, the outcomes of the analysis are reported using four tables. Figure 14 refers to the *aptitude*

dimension. All the tables share the same format. As above, the columns identify the five transitions that characterize the morphism under analysis. The six rows correspond to as many bins of the probability distribution, which coincide with the six activation levels of the hourglass model. In each figure, table (a) provides

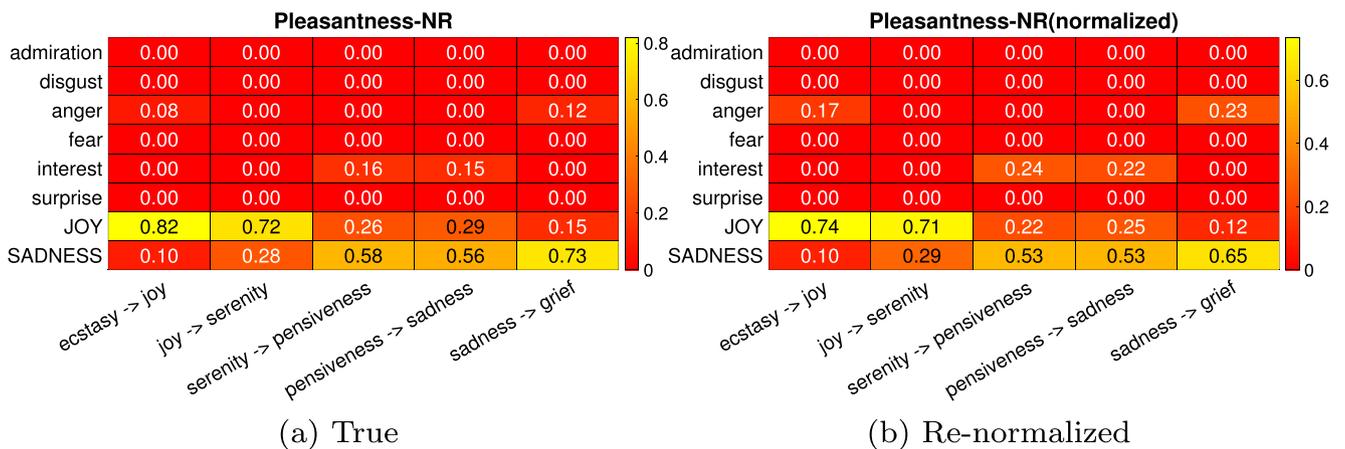


Fig. 18 a, b Pleasantness visualization using tag descriptors with non-repeated concepts

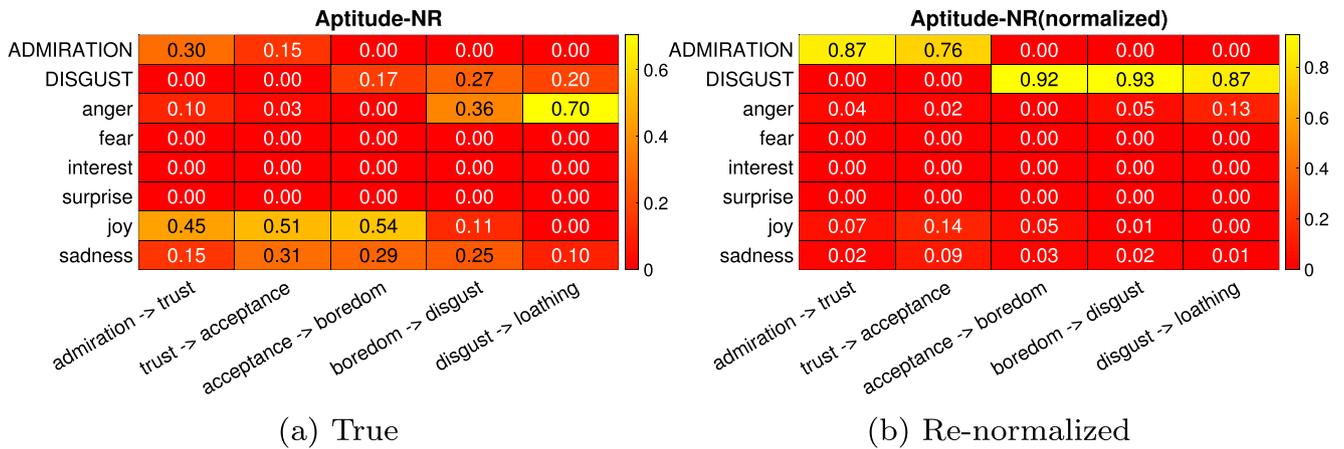


Fig. 19 a, b *Aptitude* visualization using tag descriptors with non-repeated concepts

the probability distribution of the activation levels that characterize *aptitude*. Similarly, tables (b), (c), and (d) provide, respectively, the probability distribution of the activation levels that characterize *attention*, *pleasantness*, and *sensitivity*. For clarity, the reader can refer to the schema presented in Fig. 8. Each column of a table contains the probability distribution of the activation levels computed on the concepts encountered by that path. Again, a color scale is used in all tables to facilitate a visual investigation.

Ideally, where there is coherence between the embedding space and the reference hourglass model, the four tables should confirm that the four dimensions of the cognitive model are independent. Thus, when the table concerns activation levels that do not stem from the dimension under analysis, the probability distribution should not change significantly when moving through the morphism, i.e., through the columns of the table. In Fig. 14, since *aptitude* is the dimension under analysis, tables (b), (c), and (d) are expected to show such behavior. Conversely, table

(a) is expected to show that, when moving through the morphism, the distribution of the activation levels changes accordingly. This is because the activation levels arise from the dimension under analysis. In principle, one would expect all the cells along the diagonal to be yellow, while the rest of the table should be red. The four tables in Fig. 14 indicate that, when inspected along the *aptitude* dimension, AffectiveSpace was almost coherent with the reference cognitive model. Despite a significant component of noise, table (a) mostly adheres to the ideal behavior. Tables (b), (c), and (d) show that the probability distributions involved are almost independent of the path, with only minor fluctuations observed.

Figures 15, 16, and 17 refer to the remaining three dimensions. Overall, the tables qualitatively match the expected behavior. Notably, in Fig. 17, table (a), the distribution is spread over a wide range of bins. However, the distribution remains almost constant over the columns. That is, when moving through the *sensitivity* dimension, the activation levels of *sensitivity* do not change significantly.

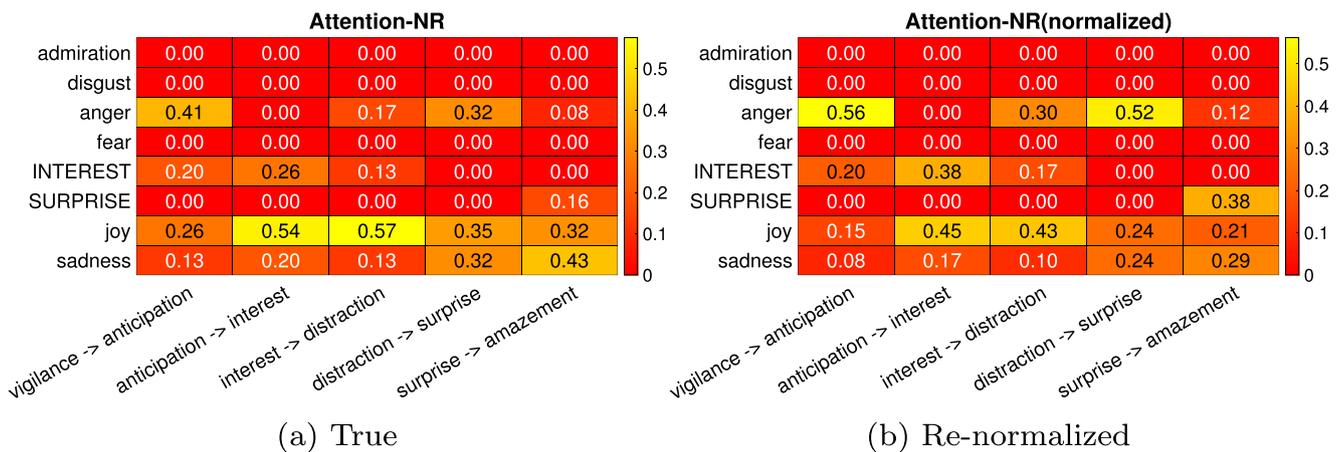


Fig. 20 a, b *Attention* visualization using tag descriptors with non-repeated concepts

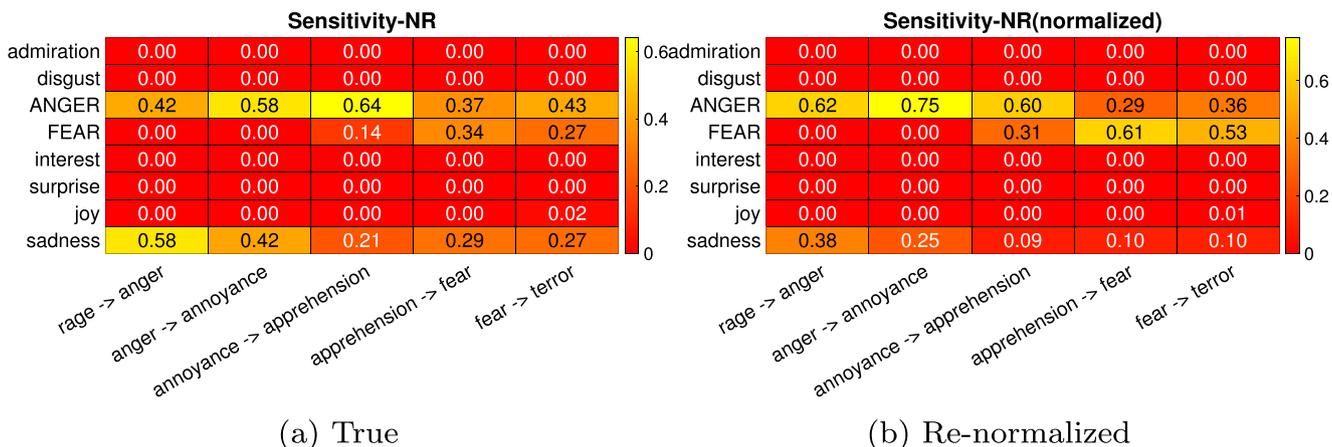


Fig. 21 a, b Sensitivity visualization using tag descriptors with non-repeated concepts

As a result, the coherence between AffectiveSpace and the hourglass model is also confirmed when focusing on these three dimensions.

Robustness Analysis

This section proposes an analysis of the robustness of the proposed approach with respect to the number of prototypes used to describe the paths. The proposed set of descriptors relies on statistical properties of the path retrieved by the RKM algorithm. In principle, the number of prototypes does not change the structure of the path. However, the procedure obtains the cognitive descriptors by projecting the prototypes onto the nearest known concept. It is therefore possible that, when the paths describe underpopulated areas of the space, more prototypes are projected onto a single known concept.

Figure 9 provides an example. The central part of the path (red crosses) goes through a part of the space that does not contain real concepts (blue dots). Following the proposed procedure, all the prototypes in this part of the 2D space are projected onto the nearest real concepts. As a consequence, multiple crosses are assigned to the same dot. Furthermore, the relative distances between prototypes and real concepts are greater than those inside densely populated regions of the space. As a result, a few concepts strongly influence the statistical properties of the path. This mechanism becomes more relevant as the number of prototypes increases.

In this section, we propose a slightly modified version of the experimental setup described in the “Tag Analysis” section. Here, the algorithm can select a real concept only once. If two prototypes are projected onto a single known concept, the information about one of the two prototypes is discarded. From a cognitive point of view, this approximation is too strong. However, this configuration removes a possible bias of the results.

The results are depicted in Figs. 18, 19, 20 and 21. Each figure has the configuration presented in the “Tag Analysis” section. The outcome of the proposed experiment is that the proposed procedure, as expected, adds noise to the description of the space. However, the trends observed in the “Tag Analysis” section are confirmed in all the presented plots.

As a result, this experimental setup empirically demonstrates that the analysis of the space is not affected by the number of prototypes used to represent the path.

Conclusion

In this paper, we presented a methodology that uses principal paths in data space and psychologically motivated descriptors in order to characterize high-dimensional spaces induced by embeddings for concepts. We focused on AffectiveSpace, visualizing a set of cognitive descriptors, which highlighted congruences between concept distributions in that space and the hourglass model of sentiment. Two complementary sets of descriptors “tag” and “distribution” were proposed. The first set models the mechanism by which human brains immediately perceive emotions, i.e., concepts characterized by the strongest component. Conversely, the “distribution” descriptors provide complete information about the sentic content of a concept. In the future, we aim to extend the analysis to more conceptual embeddings, providing a comparative study of the available embeddings.

Compliance with Ethical Standards

Conflict of interest The authors declare that they have no conflict of interest.

Informed Consent Informed consent was not required as no human or animal subjects were involved.

Human and Animal Rights This article does not contain any studies with human or animal subjects performed by any of the authors.

References

- Cambria E, Poria S, Gelbukh A, Thelwall M. Sentiment analysis is a big suitcase. *IEEE Intell Syst.* 2017;32(6):74–80.
- Li Y, Pan Q, Yang T, Wang S, Tang J, Cambria E. Learning word representations for sentiment analysis. *Cogn Comput.* 2017;9(6):843–851.
- Ofek N, Poria S, Rokach L, Cambria E, Hussain A, Shabtai A. Unsupervised commonsense knowledge enrichment for domain-specific sentiment analysis. *Cogn Comput.* 2016;8(3):467–477.
- Ma Y, Peng H, Khan T, Cambria E, Hussain A. Sentic lstm: a hybrid network for targeted aspect-based sentiment analysis. *Cogn Comput.* 2018;10(4):639–650.
- Yang H-C, Lee C-H, Wu C-Y. Sentiment discovery of social messages using self-organizing maps. *Cogn Comput.* 2018;10(6):1152–1166.
- Peng H, Cambria E, Hussain A. A review of sentiment analysis research in chinese language. *Cogn Comput.* 2017;9(4):423–435.
- Bengio Y, Ducharme R, Vincent P, Jauvin C. A neural probabilistic language model. *J Mach Learn Res.* 2003;3(Feb):1137–1155.
- Collobert R, Weston J. A unified architecture for natural language processing: Deep neural networks with multitask learning. In: *Proceedings of the 25th International Conference on Machine Learning.* ACM; 2008. p. 160–167.
- Huang EH, Socher R, Manning CD, Ng AY. Improving word representations via global context and multiple word prototypes. In: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1, Association for Computational Linguistics;* 2012. p. 873–882.
- Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector space. arXiv:1301.3781.
- Mnih A, Hinton G. Three new graphical models for statistical language modelling. In: *Proceedings of the 24th International Conference on Machine Learning.* ACM; 2007. p. 641–648.
- Tang J, Qu M, Mei Q. Pte: Predictive text embedding through large-scale heterogeneous text networks. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* ACM; 2015. p. 1165–1174.
- Wang S, Tang J, Aggarwal C, Liu H. Linked document embedding for classification. In: *Proceedings of the 25th ACM International Conference on Information and Knowledge Management.* ACM; 2016. p. 115–124.
- Ma Y, Peng H, Cambria E. Targeted aspect-based sentiment analysis via embedding commonsense knowledge into an attentive LSTM. In: *AAAI;* 2018. p. 5876–5883.
- Pennington J, Socher R, Manning CD. Glove: Global vectors for word representation. In: *Empirical Methods in Natural Language Processing (EMNLP);* 2014. p. 1532–1543. <http://www.aclweb.org/anthology/D14-1162>.
- Wilson T, Wiebe J, Hoffmann P. Recognizing contextual polarity in phrase-level sentiment analysis. In: *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing, Association for Computational Linguistics;* 2005. p. 347–354.
- Mohammad SM, Turney PD. Crowdsourcing a word–emotion association lexicon. *Comput Intell.* 2013;29(3):436–465.
- Cambria E, Poria S, Hazarika D, Kwok K. SenticNet 5: Discovering conceptual primitives for sentiment analysis by means of context embeddings. In: *AAAI;* 2018. p. 1795–1802.
- Li X, Xie H, Chen L, Wang J, Deng X. News impact on stock price return via sentiment analysis. *Knowl-Based Syst.* 2014;69:14–23.
- Cambria E, Fu J, Bisio F, Poria S. Affectivespace 2: Enabling affective intuition for concept-level sentiment analysis. In: *AAAI;* 2015. p. 508–514.
- Carlsson G. Topology and data. *Bull Am Math Soc.* 2009;46(2):255–308.
- Pearson K. Liii. on lines and planes of closest fit to systems of points in space. *Lond Edinb Dublin Philos Mag J Sci.* 1901;2(11):559–572.
- Schölkopf B, Smola A, Müller K-R. Kernel principal component analysis. In: *International Conference on Artificial Neural Networks.* Springer; 1997. p. 583–588.
- Roweis ST, Saul LK. Nonlinear dimensionality reduction by locally linear embedding. *Science.* 2000;290(5500):2323–2326.
- Kruskal JB. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika.* 1964;29(1):1–27.
- Maaten Lvd, Hinton G. Visualizing data using t-sne. *J Mach Learn Res.* 2008;9(Nov):2579–2605.
- Liu S, Maljovec D, Wang B, Bremer P-T, Pascucci V. Visualizing high-dimensional data: Advances in the past decade. *IEEE Trans Vis Comput Graph.* 2017;23(3):1249–1268.
- Ragusa E, Gastaldo P, Zunino R, Cambria E. Learning with similarity functions: a tensor-based framework. *Cogn Comput.* 2019;11(1):31–49.
- Peng X, Selvachandran G. Pythagorean fuzzy set: state of the art and future directions. *Artif Intell Rev.* 2017:1–55.
- Ferrarotti MJ, Rocchia W, Decherchi S. Finding principal paths in data space. *IEEE Transactions on Neural Networks and Learning Systems.* 2018:1–14. <https://10.1109/TNNLS.2018.2884792>.
- Hastie T, Stuetzle W. Principal curves. *J Am Stat Assoc.* 1989;84(406):502–516.
- Plutchik R. The nature of emotions: Human emotions have deep evolutionary roots, a fact that may explain their complexity and provide tools for clinical practice. *Am Sci.* 2001;89(4):344–350.
- Cambria E, Livingstone A, Hussain A. The hourglass of emotions. In: *Cognitive Behavioural Systems.* Springer; 2012. p. 144–157.
- Liu H, Singh P. Conceptnet—a practical commonsense reasoning tool-kit. *BT Technol J.* 2004;22(4):211–226.
- Strapparava C, Valitutti A, et al. Wordnet affect: an affective extension of wordnet. In: *Lrec, Vol. 4, Citeseer;* 2004. p. 1083–1086.
- Cambria E, Poria S, Bajpai R, Schuller B. SenticNet 4: A semantic resource for sentiment analysis based on conceptual primitives. In: *COLING;* 2016. p. 2666–2677.
- Cambria E, Hussain A. Sentic computing: a Common-Sense-Based framework for Concept-Level sentiment analysis. Cham: Springer; 2015.
- Bottou L, Bengio Y. Convergence properties of the k-means algorithms. In: *Advances in Neural Information Processing Systems;* 1995. p. 585–592.
- Belkin M, Niyogi P. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput.* 2003;15(6):1373–1396.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Affiliations

Edoardo Ragusa¹  · Paolo Gastaldo¹ · Rodolfo Zunino¹ · Marco Jacopo Ferrarotti² · Walter Rocchia² · Sergio Decherchi³

Paolo Gastaldo
paolo.gastaldo@unige.it

Rodolfo Zunino
rodolfo.zunino@unige.it

Marco Jacopo Ferrarotti
marco.ferrarotti@iit.it

Walter Rocchia
walter.rocchia@iit.it

Sergio Decherchi
sergio.decherchi@iit.it

- ¹ Department of Electrical, Electronic and Telecommunications Engineering, and Naval Architecture, DITEN, University of Genoa, Genoa, Italy
- ² ConceptLab, Fondazione Istituto Italiano di Tecnologia, Genoa, Italy
- ³ Computational and Chemical Biology, Fondazione Istituto Italiano di Tecnologia, Genoa, Italy