



UNIVERSITÀ DI PARMA

ARCHIVIO DELLA RICERCA

University of Parma Research Repository

Continual representation learning for node classification in power-law graphs

This is the peer reviewed version of the following article:

Original

Continual representation learning for node classification in power-law graphs / Lombardo, G.; Poggi, A.; Tomaiuolo, M.. - In: FUTURE GENERATION COMPUTER SYSTEMS. - ISSN 0167-739X. - 128:(2022), pp. 420-428. [10.1016/j.future.2021.10.011]

Availability:

This version is available at: 11381/2904048 since: 2022-01-20T14:56:50Z

Publisher:

Elsevier B.V.

Published

DOI:10.1016/j.future.2021.10.011

Terms of use:

Anyone can freely access the full text of works made available as "Open Access". Works made available

Publisher copyright

note finali coverpage

(Article begins on next page)

22 December 2024

Continual Representation Learning for Node Classification in Power-Law Graphs

Gianfranco Lombardo^a, Agostino Poggi^a, Michele Tomaiuolo^a

^a*Department of Engineering and Architecture, University of Parma, Italy*

Abstract

The recent advent of node embedding techniques enabled a more efficient application of machine learning techniques on graphs. These techniques allow each node of a network to be encoded into an arbitrary low-dimensional vector representation, which can be exploited by statistical learning models. However, the main limitation of these approaches is that the embedding task is solved as an optimization problem on a static snapshot of the graph. In a real scenario, temporal dynamics should be considered with some consequences: new nodes might join the network and get a representation of only these new ones. As a consequence, a new training step over the entire graph is required. Even more, training models with static approaches can have resource-intensive requirements, especially when dealing with large networks. In light of this, a continual feature learning that builds on top of previously already learned knowledge (previous partial embedding of the network) and well-known properties can be a solution to address both limitations efficiently in real scenarios. Our approach is suitable for graphs whose degree distribution is described by a power-law function that is a common property of real systems. This research work presents three main scientific contributions: (a) a continual feature learning meta-algorithm for node embedding, which exploits properties of power-law distribution and spaces alignment techniques; It is suitable with any traditional node embedding techniques that relies on embedding spaces (b) we demonstrate empirically, by performing node labeling tasks, that a lightweight solution to encode new nodes, based on limited knowledge of the embedding of the network hub-nodes, can provide comparable or better performances, with respect to static approaches. (c) Finally, we experimented our algorithm in the temporal graphs domain and we achieved better results in node classification compared with other state of the art techniques.

Keywords: Embedding, Node Embedding, Incremental learning, Continual Learning, Dynamic Networks, Representation learning

1. Introduction

Graphs (or networks) representations are fundamental in modeling many systems, particularly for expressing relational knowledge about interacting entities. Thus, they are ubiquitous in a wide range of disciplines, such as Computer Science, Sociology, Finance, Biology, and others. Moreover, in addition to being useful as knowledge repositories, these data structures also play a key role in Knowledge Discovery within the application of Machine Learning techniques. Several examples are available in the scientific literature of different research fields: in [1] the authors modeled the interaction between proteins as a network, to predict a correct label for each node, for describing its functionalities; in [2] the authors use a temporal network describing the US stock market to cluster correlated stocks' behaviors, in order to predict financial crises; in [3] an online community of patients is exploited to extract latent information about the emotional effects of their pathology. The main issue when applying machine learning on graphs is finding a way to encode structural information in the feature space required by learning models.

Representation learning is the branch of machine learning that aims to solve this problem by generalizing low-dimensional representations learning for data structures or sequences (e.g. text, graph). State of the art is learning these embeddings by solving an optimization problem that seeks to preserve local and global properties [4]. The first attempt in this field is the Neural language model [5] and its most famous breakthrough implementation for word embedding: Word2Vec [6]. In this model, a shallow neural network is trained to predict the missing word given the context or the contrary of predicting the context given the word. The training task is an auxiliary task to get feature vectors associated with each word, named representation learning. In particular, Word2vec exploits the Skip-gram model to optimize the objective function using Stochastic Gradient Descent (SGD) as an iterative learning function and hierarchical Softmax. Both are used in the form of back-propagation on a single hidden layer feed-forward neural network with a linear activation function. The Skip-gram model has been adapted in the last years for node embedding in networks, by using streams of random walks instead of text corpora.

34 Although these approaches raised state-of-art benchmarks in both fields,
35 they present a limitation when dealing with dynamical contexts. They are
36 designed to work on a static context, for example, a snapshot of a graph at
37 a certain time. In a real scenario, a new node might join the network or a
38 word in a text might be unseen at training time. It is necessary to run the
39 entire representation learning process again, to get a feature vector for these
40 new elements. Addressing this issue in a dynamic context is still challenging
41 and of increasing interest in Machine Learning, as demonstrated by the high
42 number of papers published in the last two years.

43 In this research work, we propose a meta-algorithm that aims to learn
44 features of new nodes incrementally joining a graph, by using a continual
45 learning approach, partial embedding alignment and exploiting properties of
46 power-law graphs. This approach is suitable for graphs whose degree distri-
47 bution follows a power-law function. However, several works [7] have primar-
48 ily demonstrated that real systems tend to present this distribution in their
49 node degree. We demonstrate that it is possible to embed new data by learn-
50 ing features incrementally, building on top of previously learned knowledge
51 and well-known properties, with comparable or even outperforming results
52 with respect to the static methods, without training the whole model again.
53 Moreover, we experimented our algorithm for Node Classification on tempo-
54 ral graphs achieving better results than other state of the art methods.

55 The paper is organized as follow: in Section 2, we present the current state
56 of the art in this domain, in Section 3 we introduce some properties related
57 to power-law graphs that are involved as an assumption in our algorithm.
58 In Section 4, we describe the proposed solution, named WalkHub2vec and
59 finally in Section 5 the experimental part, and the results are presented.

60 **2. State of the art**

61 The increasing interest in Representation Learning for graphs in the Ma-
62 chine Learning community is due to the idea that exploiting links among
63 data in the form of a graph can increase learning performances for several
64 tasks. In particular, the witnessed rapid growth in knowledge graph (KG)
65 construction and the application made this task crucial in several domains.
66 However, finding a way to incorporate information about the structure of a
67 network into a statistical model is still challenging because of the heterogene-
68 ity of this information. For example, in the case of link prediction, one might
69 want to focus more on structural properties to encode pairwise properties be-

70 tween nodes (e.g. common friends in a social network). On the contrary, in
71 the case of node classification, one might want to preserve more global prop-
72 erties in the node description (e.g. nodes sharing many labels need to be
73 similar). Moreover, in knowledge graphs, each edge is represented as a triple
74 of the form (head entity, relation, tail entity), also called a fact, indicating
75 that a specific relation connects two entities. Thus, the underlying symbolic
76 nature of such triples usually makes this kind of graph hard to manipulate
77 for machine learning tasks [8].

78 The idea behind the representation Learning is to learn a mapping func-
79 tion that embeds nodes as points in a low-dimensional vector space by solving
80 a downstream task. Several approaches have been presented in the scientific
81 literature, with multiple and separate lines of research [9] [10]. In knowledge
82 graphs, the key idea is to embed components of a KG including entities and
83 relations into this kind of continuous vector spaces, to learn features that
84 preserve the inherent structure of the graph. Those embeddings can further
85 be used to benefit all kinds of tasks, such as KG completion [11],[12], entity
86 resolution [13], node classification [14], [15] and link prediction [16].

87 The State of the Art is represented by several methods that learn static
88 node representations combining neural networks, random walks in graphs
89 and the neural language model Skip-gram [17]. A classic approach to node
90 embedding is DeepWalk [18], which captures second-order proximity. In this
91 algorithm, random walks are modeled as sentences and later are fed into
92 the Skip-gram model. Values from the hidden layers are the resulting node
93 embedding vectors.

94 The LINE algorithm [19] addresses the efficiency issues, that previous ap-
95 proaches have when applied on large networks, by improving scalability with
96 the use of negative sampling and asynchronous stochastic gradient descent
97 to solve the optimization problem.

98 Finally, another relevant algorithm is Node2vec [20], which is an impor-
99 tant modification of DeepWalk with significant performance gains. It intro-
100 duces a parametrization of random walks generation by combining DFS and
101 BFS-like behavior and also negative sampling. A common formulation of the
102 optimization problem addressed by these algorithms can be summarized by
103 Equation 1. where f is the function that maps each node $u \in V$ in the vector
104 space and $N_s(u)$ is the neighborhood of node u , sampled in different ways by

105 the algorithms.

$$\max_f \sum_{u \in V} \log Pr(N_s(u)|f(u)) \quad (1)$$

106 However, although much progress has been reached in terms of perfor-
107 mances, scalability over large networks is still problematic. It can become a
108 limit even with just thousands of nodes when the network is dense. Addi-
109 tionally, the static nature of these methods can hardly apply to the analysis
110 of real systems that are intrinsically dynamic.

111 In the last two years, several works have tried to address the issue of
112 node embedding in a dynamic context, and this increasing production of
113 different solutions highlights the importance of the topic both in academia
114 and industry. The most used approach in the literature is based on modeling
115 dynamics as a sequence of static snapshots.

116 The first algorithm, based on network snapshots, is the Continuous-time
117 Dynamic Network Embedding (**CTDNE**) [21], proposed for the link predic-
118 tion task. It relies on the concept of temporal random walks, using times-
119 tamps to preserve the time order of edges and thus their temporal dependen-
120 cies. It can preserve temporal structural properties for link prediction, but
121 it requires retraining at every change in the network.

122 Other successful snapshots-based works, relying on a different approach,
123 are **DynGEM** [22], the one proposed in [23], **Dynnode2vec** [24] and **Dyn-**
124 **graph2vec** [25]. Dyngem and Dyngraph2vec estimate node representation
125 with a deep autoencoder and an LSTM, respectively, by initializing net-
126 work weights with the previously computed embedding at time $t - 1$. Both
127 algorithms address the problem of graph visualization and link prediction
128 without node classification. Their main limitation is their use of their own
129 learning model (PropSize), so they only work with this embedding method.
130 The algorithm proposed in [23] is an extension of the Skip-gram model and,
131 in particular of the LINE algorithm. Its interesting contribution is that when
132 a new node joins the network, changes in the representations of other nodes
133 are limited. In light of this, it computes a new representation with new walks
134 only for a selected set of vertices. It offers impressive performance in several
135 tasks, including a simple task of node classification, independently to the
136 static algorithm used for random walk generation.

137 Finally, Dynnode2vec can be seen as a combination that takes the advan-
138 tages of CTDNE and the algorithm in [23]. It uses the concept of evolving
139 walk generation and computes new representation only for new nodes that

140 join the graph in a particular timeframe by initializing the hidden weights of
141 the neural network with the ones used to embed the snapshot at time $t - 1$.

142 More recently, in [26], the authors proposed an algorithm, named Credit
143 probing network embedding, that aims to compute node embedding for evol-
144 ving networks with a partial monitoring of the network. The basic idea is to
145 update each node vector dealing with it as a timestamps vector and to exploit
146 a probing system to decide which node and when it should be updated. How-
147 ever, in this case, the evaluation task is represented by the link prediction
148 task, and the source code is not currently available.

149 Finally, in [27], the authors propose **tNodeEmbed**, an algorithm based
150 on Node2Vec that learns new nodes representation by aligning snapshots us-
151 ing Orthogonal Procrustes [28] and then training an LSTM for a given task
152 to optimize feature learning. It is not easy to assert which one is the best,
153 because they are experimented on different datasets and with some gaps for
154 what concerns other more difficult tasks like Node classification. The only
155 tNodeEmbed takes into consideration the task of multi-class node classifica-
156 tion, i.e., the benchmark task allowing a comparison with static approaches.
157 This last algorithm reaches very interesting performance, but a complete
158 comparison with the static case remains difficult. The authors experimented
159 with tNodeEmbed on different datasets and the used parameters for embed-
160 ding are unclear.

161 From direct experiments and the analysis of recent literature, we can as-
162 sert that the idea of dynamic network embedding based on snapshot does
163 not solve scalability issues and retraining is time-consuming, although less
164 than the static case under certain conditions. In particular, these are im-
165 portant limitations when dealing with real networks that can often be large
166 and dense. We have to be able to maintain previous embeddings and some-
167 times the set of previous walks whenever there is a change. It is possible
168 to conclude that several issues and open problems are present in dynamic
169 network embedding. However, considering the recent growth of this research
170 line, all of the works reviewed in this section propose interesting approaches
171 with important progress that trace a way to more generalized solutions.

172 In this work, we take into account specifically the multi-class node labeling
173 task in light of *(i)* the reported lacks in state of the art, and *(ii)* the different
174 issues which this case presents, with respect to the more studied case of link
175 prediction.

176 3. Properties of power-law graphs

177 WalkHub2Vec finds its theoretical justifications in several properties of
178 power-law graphs (Scale-free networks). The importance of this particular
179 category of networks has been widely discussed in several research fields (E.g.
180 [29],[30],[31]). The main reason is related to the fact that it fits the relation-
181 ships among entities in many real systems, in various domains. The difference
182 between a random network and a power law one lies in the probability dis-
183 tribution, which better describes a node’s choice with a particular degree.
184 In random networks, the randomness of links generation brings to a Poisson
185 distribution and a lack of a component composed of nodes with a higher
186 degree than other nodes, i.e., the hub nodes [7]. On the other hand, real
187 systems are often better described by a power-law because only a few nodes
188 present a high degree, while most of the remaining nodes have a small degree.
189 When a system presents a power-law degree distribution, it means that the
190 probability of randomly selecting a node with degree K is well approximated
191 by Equation 2

$$P_k \sim k^{-\gamma} \quad (2)$$

192 This difference has an important consequence on how networks tend to
193 grow. In power-law graphs, the evolution is well described by the Barabási-
194 Albert Model [32]. This model assumes that in power-law graphs (or Scale-
195 free graphs) the growth is described by a property called **Preferential At-**
196 **tachment**. This property claims that the probability $\Pi(k)$ of a new node
197 to be connected with node i depends on the degree k_i , as in Equation 3.
198 That means that the growth in this kind of graph follows a probabilistic
199 mechanism where a new node is free to connect to any node in the network,
200 whether it is a hub or has a single link.

201 However, if a new node has a choice between a degree-two and a degree-
202 four node, it is twice as likely that it connects to the degree-four node [7].
203 This behavior related to the power-law distribution is well known also as the
204 Pareto principle, or 80/20 rule [33]. This rule can be summarized, saying
205 that roughly 80% of the effects come from 20% of the causes. For exam-
206 ple, as Pareto noticed in the 19th century, 20% of the population earned
207 most of the money, while most of the population (80%) earned rather small
208 amounts. Thus, it is likely that new incomes are destined to that 20% of
209 people. Similarly, if a new node joins a network, it will likely be connected
210 with the graph’s 20% of hub nodes. This property relies on the long-tail of
211 the power law distribution: in this kind of graphs, it is possible to observe

212 that most nodes have a small degree while few nodes have a big degree (the
 213 so called long-tail). For this reason, as illustrated in [7] in power-law graphs
 214 is present a giant component that is composed by hub nodes (nodes with a
 215 high degree) that connect most of the normal nodes with a small degree.

$$\prod(k_i) = \frac{k_i}{\sum_j k_j} \quad (3)$$

216 4. WalkHub2Vec

217 In this section we describe the incremental feature learning algorithm for
 218 node embedding. In light of the properties summarized in Section 3, we take
 219 into account power-law graphs. First of all because WalkHub2Vec aims to be
 220 a solution to embed real systems graphs, thus this choice does not represent
 221 a limitation. Secondly, to the best of our knowledge, these properties have
 222 not been taken into account by previous works. Finally, as we demonstrate
 223 in the result section, these properties represent a useful tool to address the
 224 dynamical case of a new node joining the network and in need for a vector
 225 representation. Also, we believe that an important part of the information
 226 required by the node classification can be directly found in the hubs compo-
 227 nent of each graph. The algorithm that we present is based on DeepWalk for
 228 simplicity, but it can be used with other embedding solution (e.g. Node2Vec
 229 or LINE) and also combining them in different moment. Additionally, being
 230 it based on a static embedding of a network, it is possible to retrieve past
 231 trained embeddings of a network to obtain the representation of new nodes
 232 without retraining. Every time a new node has to be embedded with respect
 233 to the previous static embedding, we seek to optimize a variant of the problem
 234 in Equation 1 specifically for incremental feature learning. This optimization
 235 maximizes the probability of observing a second-order neighborhood among
 236 the hubs component for the new node i , conditioned on an aligned feature
 237 representation between the target original space and a lightweight embedding
 238 of the hubs, or in other words the embedding of node i with respect to the
 239 original embedding space:

$$\max Pr(N_H(i)|R \cdot f_H(i)) \quad (4)$$

240 Where R is the rotational matrix necessary to align the drifted embedding
 241 space and is later discussed in detail. Being A defined as the static algorithm

242 chosen for the unique snapshot of the graph, WalkHub2Vec has the following
 243 steps:

- 244 1. It computes the static embedding of the network G with A following the
 245 specific parameters of the considered algorithm and solving the general
 246 optimization problem in Equation 1. Given a graph G we define also
 247 $E_{G\{G\}}$ as the embedding of G made with respect to G itself.
- 248 2. The algorithm extracts the hub component of the network by selecting
 249 the sub-graph H induced by the percentage λ of nodes with the highest
 250 degree. The parameter λ is set as default equal to 20 in light of the
 251 Pareto principle. For construction $E_{H\{G\}} \subset E_{G\{G\}}$ without any other
 252 computation (Figure 1a).
- 253 3. When a new node i joins the network, it is considered a graph composed
 254 by the sub-graph H plus the node i : $H + i$. In the rare case when the
 255 new node in the power-law graph has not links with the hubs, a single
 256 random walk from i to one of the hubs is considered in the construction
 257 of $H + i$. Given the smaller network $H + i$ a lightweight embedding
 258 $E_{H\{H+i\}}$ is computed with A or an arbitrary algorithm (Figure 1b).
- 259 4. In order to get the embedding of node i with respect to the original
 260 embedding space, the algorithm proceed to align $E_{H\{H+i\}}$ with $E_{H\{G\}}$
 261 by learning the optimal rotation matrix. The alignment involves only
 262 nodes' representations who belong to the second-order neighborhood of
 263 i in the hub component (Figure 1c). Alignment details are discussed
 264 later in 4.1.
- 265 5. Once the rotational matrix R is computed, the embedding of i with
 266 respect of the original embedding space is the dot product between R
 267 and the embedding of i with respect to $H + i$ space (Figure 1d).

268 4.1. The orthogonal Procrustes problem

269 The orthogonal Procrustes problem [28] is a matrix approximation tech-
 270 nique that aims to learn the orthogonal rotational matrix R which closely
 271 maps a matrix A to a matrix B (Equation 5). It is based on the Frobenius
 272 norm and a closed-form solution is provided with the Single Value Decom-
 273 position (SVD).

$$R = \operatorname{argmin}_{\Omega} \|\Omega A - B\|_F \text{ subject to } \Omega \Omega^T \quad (5)$$

274 In WalkHubs2Vec, matrix A is represented by the embedding of the second-
 275 order neighborhood of node i in $E_{H\{H+i\}}$ and respectively the correspondents

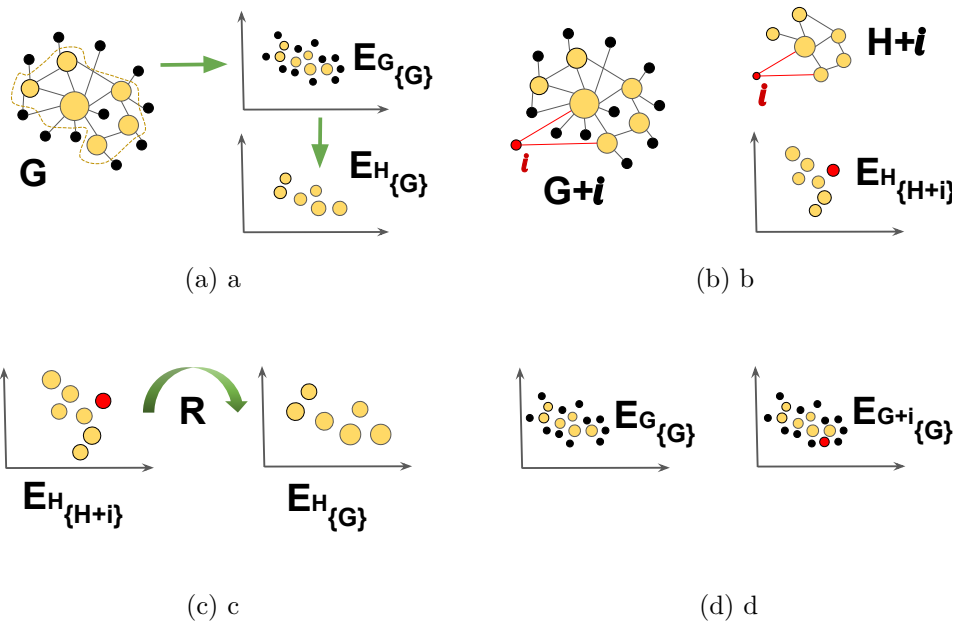


Figure 1: (a) First and second steps of the algorithm; The embedding of the hub component is extracted as a subset of the entire embedding E_G (b) Third step: A new node i joins the network and a new embedding of only the H sub-graph is computed; (c) The two embedding spaces have all nodes in common except for the new node i . They have to be aligned to get i coordinates with respect to the original space. Embedding alignments by learning the optimal rotation matrix; (d) Node i with its resulting embedding in the original space

276 in $E_{H_{\{G\}}}$. In order to increase the quality of the alignment the two matrix are
277 previously pre-processed following the full Procrustes superimposition [34]. It
278 requires that matrices, that can be seen as geometric shapes to be aligned,
279 are previously uniformly scaled and translated to have the average value in
280 the origin.

281 5. Results

282 In this section we provide an overview of the datasets and methods which
283 we will use in our experiments. Code and data to reproduce our results are
284 available on GitHub (see footnote 1). In order to evaluate WalkHubs2Vec
285 and the related methodology, we performed three different kind of node la-
286 labelling tasks: multi-label classification, multi-class classification and node
287 classification with a dynamic graph. The difference between the first and the
288 second is that in the first one each node can have more than one label at
289 the same time and each one has to be correctly predicted. For both tasks,
290 with the aim of a fair comparison between our method and the static tech-
291 niques, we used the exact experimental procedure as in [35, 36], DeepWalk
292 [18] and Node2Vec [20] and their common datasets. Finally, we evaluated
293 WalkHubs2Vec in a dynamic context with temporal edges and nodes, making
294 a comparison with CTDNE[21], tNodeEmbed [27] and DySAT [37]. For each
295 task and experiment, we randomly selected a percentage of the labeled nodes
296 between 10 and 90%, and use them as training data. The rest of the nodes
297 are used as test set Specifically, we repeated 10 times, and we report the
298 average performance in terms of both Macro-F1 and Micro-F1. The machine
299 learning model used for all the tasks is a One-Vs-Rest logistic regression. ¹

300 5.1. Datasets

301 For the multi-label classification we considered three networks: PPI and
302 Wikipedia. We selected these networks as a benchmark because have also
303 been used to evaluate Node2Vec and DeepWalk in their respective papers.

- 304 • **PPI**: Protein-Protein interactions for Homo Sapiens, is a network with
305 3,890 nodes, 76,584 edges and 50 different labels, where labels represent
306 the biological states.

¹<https://github.com/gfl-datascience/walkHub2vec>

- 307 • **Wikipedia:** This network is composed by 4,777 nodes with 40 differ-
 308 ent labels and 184,812 edges. It is a co-occurrence network of word
 309 appearing in the first million bytes of the Wikipedia dump. Labels are
 310 Part-of-Speech (POS) tags associated to each word.
- 311 • **Blogcatalog:** It is a social network of relationships among blogger
 312 authors on the BlogCatalog website. The network has 10,312 nodes,
 313 333,983 edges and 39 different labels that represent blogger interest
 314 inferred through metadata.

315 For the multi-class classification task we used Cora, that is one of the most
 316 famous in literature for this task (E.g. [38],[39], [27]) with its various ver-
 317 sions. It is a citation network composed by papers about different research
 318 areas of Computer Science and their citations. We used the largest version
 319 available with all of the nodes and the labels [40]. The network has 23,567
 320 nodes, 93965 edges and 10 classes, which represent the topic of the articles.
 321 We exploited Cora also for the node classification task in dynamic graphs.
 322 We used temporal directed edges that represent citations from one paper to
 323 another, with timestamps of the citing paper’s publication date as in [27]. We
 324 trained a starting embedding with papers between 1900 and 1990 and then
 325 we used papers until 1999 as new nodes and edges coming-up in the network.
 326 We used the same methodology also to evaluate CTDNE, tNodeEmbed and
 327 DySAT.

328 5.2. Experimental setup

329 In light of the reasons explained in Section 2 about a necessary fair com-
 330 parison, we used the same datasets used to evaluate Node2vec and DeepWalk
 331 in their respective papers. These datasets have not temporal information but
 332 this has not been a significant limitation for our evaluation, since we do not
 333 aim to solve only a temporal problem but also an incremental one, when a
 334 new generic node joins the network and when the entire graph is unknown
 335 at training time. Under this hypothesis, for the multi-label and multi-class
 336 tasks, we randomly selected the 10% of the nodes in each dataset to simulate
 337 their appearance (I set). These nodes are sampled by a non-uniform distri-
 338 bution over the degree that makes nodes with smaller degree more likely to
 339 be sampled. This choice is also motivated by the idea that the leaf nodes
 340 with few links are often under evaluated by the sampling edges techniques
 341 used in the Skip-gram model. However, also Hub nodes can be sampled in

342 the (I set) although is less frequent that a node joins a network directly
343 with an high degree. In order to prove the effectiveness of our methodology,
344 we also treated the moments at which they join the network as totally in-
345 dependent events. In other words, we do not consider the possibly existing
346 edges in the sub-graph induced by the I set. The rest of the network rep-
347 represents the first static snapshot used as input for WalkHubs2Vec. As base
348 for the meta-algorithm we used predominantly DeepWalk and in one case
349 also Node2vec as demonstration that our algorithm and the methodology
350 can deal with any node embedding technique. We also made a compari-
351 son with these two algorithms without the use of the meta-algorithm and
352 when the graph is entirely processed with these algorithms. We used the
353 same hyper-parameters proposed in the respective papers of DeepWalk and
354 Node2Vec, without any optimization of these parameters. Each experiment
355 is composed by 10 independent runs where each time nodes for the different
356 sets are sampled.

- 357 • Dimension = 128
- 358 • Walk lenght = 10
- 359 • Number of walk = 80
- 360 • P and Q (Node2vec) both equal to zero to compare with DeepWalk

361 We also used in each case the parameter $\lambda = 20$, in this way the static
362 network is divided into two groups: 20% of hubs and the remaining 80%
363 unconsidered for the embedding.

364 5.3. Multi-label node classification

365 In this section we report the results of the multi-label classification on
366 PPI, Blogcatalog and Wikipedia. We experimented our algorithm with the
367 previously introduced methodology by selecting different portions of nodes
368 as training set, in the range between 10% and 90%. However, we report
369 only results for the case of 50%, because we are interested in evaluating
370 how features, computed incrementally on the 10% of the nodes of I set, can
371 affect the prediction results. This condition became more evident when we
372 randomly selected large portions of the network in the 10 runs. Results are
373 presented in Figure 2.

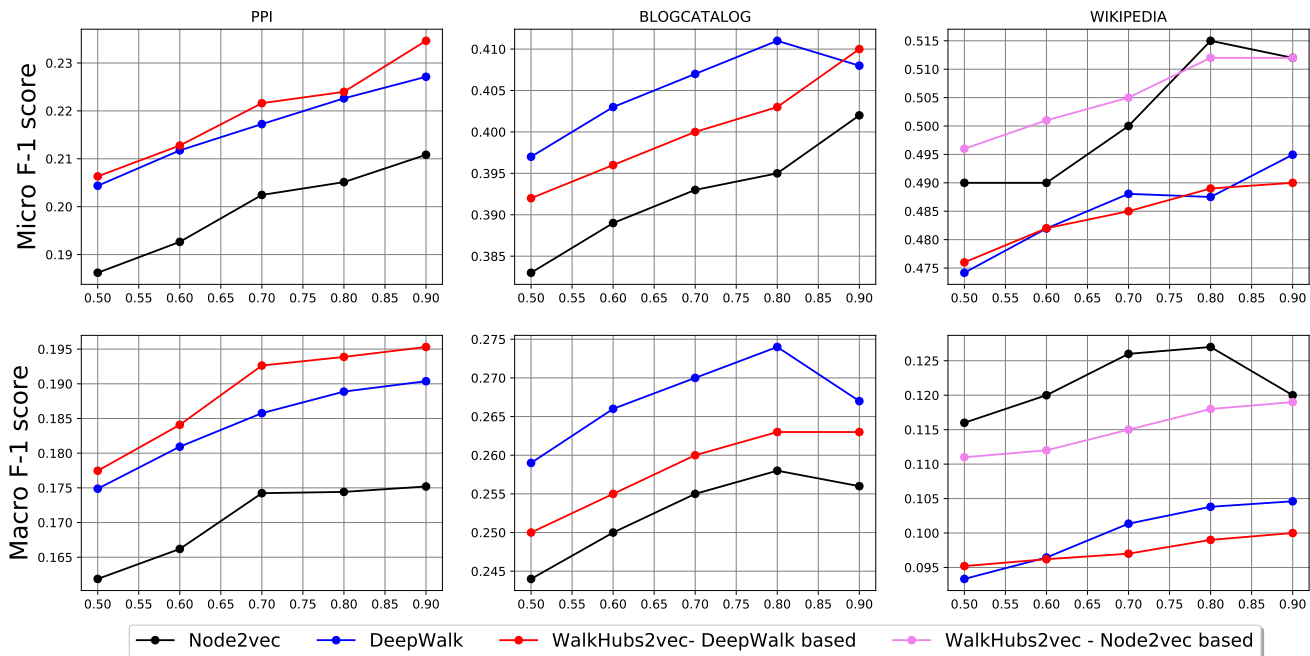


Figure 2: Performance evaluation of different benchmarks on varying amount of labeled data used for training. The x axis denotes the fraction of labeled data, whereas the y axis in the top and bottom rows denote the Micro-F1 and Macro-F1 scores respectively.

374 5.3.1. Protein-Protein Interactions

375 On this dataset, WalkHubs2Vec performs slightly better than Node2Vec
 376 and Deepwalk both with Micro and Macro F-1 score. In particular it is
 377 interesting the result when the selected percentage of training is the 90%
 378 and probably most of the nodes of I set have been considered in the 10 runs.

379 5.3.2. BlogDatalog

380 On this dataset, WalkHubs2Vec performs in a comparable way with Deep-
 381 Walk with an average divergence of -0.004 and -0.009 in terms of respectively
 382 micro and macro F-1 score.

383 5.4. Wikipedia

384 Unlike the two previous networks, on Wikipedia Node2Vec performs in
 385 general better than DeepWalk. In light of this, we present results obtained
 386 with WalkHubs2Vec based on DeepWalk but also based on Node2Vec. In

387 both cases WalkHub2Vec performs in a comparable way than the static tech-
388 niques. The average divergence between Node2vec and the incremental case
389 is 0.003 and -0.006 for micro and macro F-1 respectively. As expected, the
390 version of WalkHubs2Vec based on Node2Vec performs better than the one
391 based on DeepWalk as reflection of the static performance of the algorithms.

392 *5.5. Comparative analysis for multi-label classification*

393 Several features of the selected datasets should be considered, to get an
394 overall evaluation of our algorithm in the multi-label classification task. First
395 of all, each network presents a power-law distribution for node degree, but
396 Wikipedia graph has a more significant cut-off than the others. Practically,
397 on this dataset, the power-law distribution is more acceptable after a certain
398 degree (around 10), while before this value, a log-normal distribution is more
399 evident. In our opinion, this fact is what makes the difference with results
400 we obtained with the PPI dataset since the density of the two graphs is quite
401 similar (PPI: 0.005, Wikipedia: 0.003) and the average node degree (PPI:
402 19, Wikipedia: 15). Another critical factor is how the topological structure
403 affects the results and how the number of labels that have to be predicted.
404 We analyzed this aspect more by comparing the three networks in terms of
405 modularity for community detection and connected components. In PPI,
406 where we have better results than the other algorithms, the number of labels
407 is 50, modularity is 0.337, and it is possible to find 43 different clusters (or
408 communities) and 35 connected components. That means that nodes tend
409 to arrange in groups with the same labels, and thus, the classification task
410 results to be more simple. In the Wikipedia network, with a Modularity
411 equal to 0.2, we found only 12 clusters, a single connected component with
412 40 different labels that have to be predicted. Finally, Blogcatalog resulted
413 in a more difficult case: with a Modularity equal to 0.24 we found 7 differ-
414 ent clusters, one single connected component when the number of labels is
415 39. In light of this, we hypothesize that exploiting the power-law properties
416 can make a significant difference when the labels reflect also the topological
417 structure of the network, and on the other hand, have comparable results in
418 the other cases. However, in these last cases, having comparable results by
419 using only the 20 percent of nodes to compute the embedding is an interesting
420 result that should be more analyzed in future works.

421 *5.6. Multi-class node classification*

422 In this section, we report the results of the multi-class classification on
 423 Cora with WalkHubs2Vec based on DeepWalk with the static case where
 424 the entire network is embedded with Node2Vec and DeepWalk. In this case,
 425 WalkHubs2Vec outperforms the baseline methods both in micro and macro
 426 F-1 score. The average gain of WalkHubs2Vec is about 0.01 in both metrics.
 427 Results are presented in Figure 3.

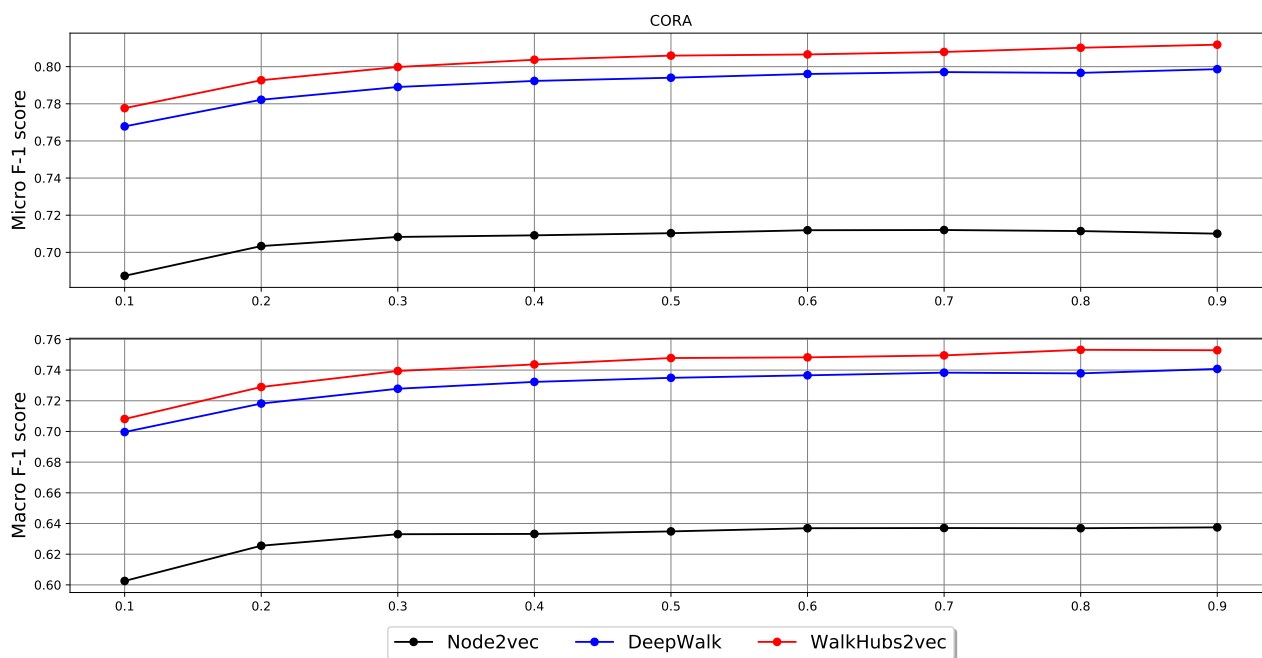


Figure 3: Performance evaluation on Cora of different benchmarks on varying the amount of labeled data used for training. The x axis denotes the fraction of labeled data, whereas the y axis in the top and bottom rows denote the Micro-F1 and Macro-F1 scores respectively.

428 *5.7. Node classification with temporal graphs*

429 WalkHubs2Vec exploits a continual learning approach since it is designed
 430 to learn a representation of nodes and parts of a graph that are not entirely
 431 known at training time. For this reason, we experimented with the algorithm
 432 in the most common scenario of temporal graphs. Indeed, in this context,
 433 the graph structure is dynamic with timestamps associated with each edge

	Micro f-1 score	Macro f-1 score
CTDNE	0.7194	0.6397
DySAT	0.4828	0.1312
tNodeEmbed	0.7740	0.7078
WalkHubs2Vec	0.8101	0.7532

Table 1: Performance results of WalkHubs2Vec vs. baselines for the node classification task over the dynamic Cora dataset.

434 that determines the moment an edge comes up to be part of the graph. As
 435 a consequence, new nodes may join the network.

436 The Cora dataset exploited in the previous section can be used also to model
 437 a temporal scenario by using as timestamps the publication date of the pa-
 438 pers, as already presented in [27]. We experimented WalkHubs2Vec on this
 439 dataset using as baselines three algorithm designed for temporal graph em-
 440 bedding: tNodeEmbedding [27], CTDNE [21] and DySAT [37]. These three
 441 algorithms have been already presented in the Literature review 2.

442 We used papers from 1900 until 1989 to create the starting graph and then
 443 the remaining ten years as new nodes and edges that join the network. In
 444 the starting graph, there are 21974 nodes and 64991 edges. Temporal items
 445 are 1593 new nodes and 28975 new edges. We computed the embedding
 446 representation for each node and then we performed a classification task to
 447 predict one of the ten possible labels. Classification has been performed us-
 448 ing the One-Vs-Rest Logistic classifier as in the previous experiments except
 449 for tNodeEmbed because it learns the embedding while dealing with a clas-
 450 sification task performed with an LSTM deep neural network. We evaluated
 451 WalkHubs2Vec and the baselines in terms of Micro-F1 and Macro-F1 scores.
 452 Table 1 shows the results we achieved on this task. WalkHubs2Vec outper-
 453 forms the baselines both in micro and macro f1-score. All the experiments
 454 have been repeated for 10 runs and the results represent the average score.

455
 456 CTDNE, DySAT and tNodeEmbedding exploits Node2Vec as node em-
 457 bedding engine for every single snapshot, for this reason, we compared our
 458 algorithm using the same embedding technique. Since Node2Vec is the basis
 459 of all the algorithms, we selected the same embedding hyper-parameters for
 460 all and they are the same of the previous sections and from [20]. We did
 461 not optimize the other parameters available within the algorithms and this
 462 can represent a reason why in particular DySAT under-performs remarkably

463 with respect of the other algorithms. DySAT and tNodeEmbed train neural
 464 network-based models while performing representation learning of nodes and
 465 for this reason, the number of epochs represents an important parameter. We
 466 decided to train CTDNE, tNodeEmbed, and DySAT for 50 epochs, because
 467 over this number we experimented with some computational issues right with
 468 DySAT. Another thing that we desire to report is that DySAT is designed
 469 to learn the distribution of edges to perform only link prediction although it
 470 computes the temporal embedding of each node.

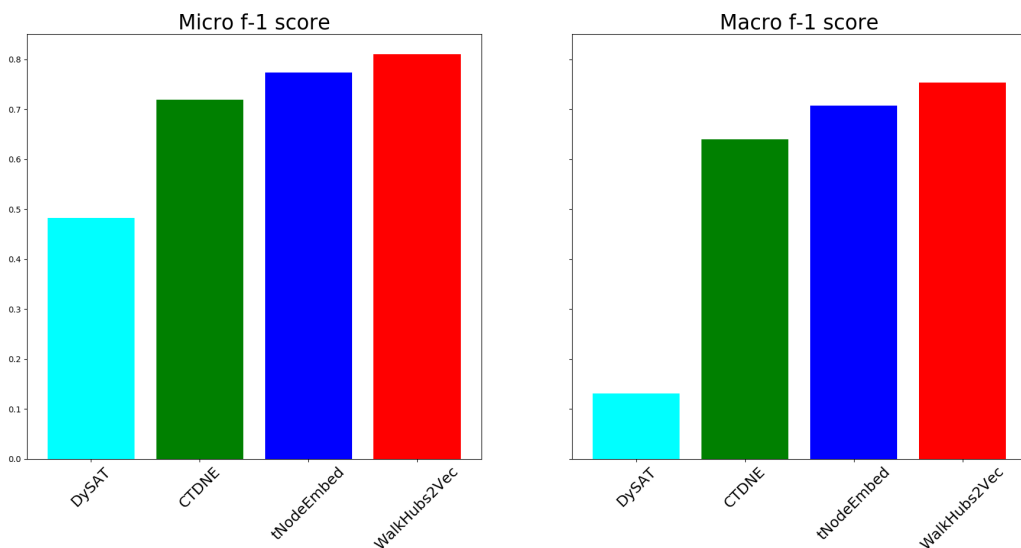


Figure 4: Performance evaluation between WalkHubs2Vec and the baselines on the temporal Cora dataset.

471 It is interesting the comparison with tNodeEmbed because a common
 472 aspect with WalkHubs2Vec is the idea of embedding alignment, that in the
 473 first one is computed over the entire graph while in our algorithm is performed
 474 considering only the giant component. Moreover, tNodeEmbed has been
 475 evaluated with its built-in deep neural network for the classification task
 476 and not with the one-vs-rest logistic classifier that is a less powerful model.
 477 This last aspect leads us to the idea that considering similarities between the
 478 two approaches and the theoretic advantage of the LSTM neural network
 479 for tNodeEmbed, one possible reason for our better results can rely exactly
 480 upon the importance of making attention to the hubs of the network while

481 learning the new nodes' representation. We plan to further investigate this
482 aspect in our future works in order to quantify this importance.

483 **6. Conclusions**

484 In this research work we have addressed the problem of continual feature
485 learning for node embedding in power-law graphs. The work is focused on the
486 idea of taking advantages of the scale-free property of this category of graphs,
487 to compute embedding of new nodes without retraining the learning model
488 in the node classification context. We propose also an implementation of this
489 methodology, WalkHubs2Vec, that can deal with any embedding techniques.
490 Results demonstrate how the combination of partial embeddings based on the
491 hubs component and embedding alignment can solve the problem with good
492 results in terms of features quality for the new nodes. WalkHubs2Vec reaches
493 equal or slightly better results when compared to well-known static methods,
494 as Node2vec and DeepWalk and better results when compared with state-
495 of-the-art techniques for dynamic graph embedding. Future developments
496 are related to performing different tasks with this methodology (e.g. link
497 prediction). Moreover, it would be interesting to use a similar approach also
498 in word embedding and other embedding contexts.

499 **References**

- 500 [1] P. Radivojac, W. T. Clark, T. R. Oron, A. M. Schoes, T. Wittkop,
501 A. Sokolov, K. Graim, C. Funk, K. Verspoor, A. Ben-Hur, et al., A large-
502 scale evaluation of computational protein function prediction, *Nature*
503 *methods* 10 (2013) 221.
- 504 [2] A. Kocheturov, M. Batsyn, P. M. Pardalos, Dynamics of cluster struc-
505 tures in a financial market network, *Physica A: Statistical Mechanics*
506 *and its Applications* 413 (2014) 523–533.
- 507 [3] G. Lombardo, P. Fornacciari, M. Mordonini, L. Sani, M. Tomaiuolo,
508 A combined approach for the analysis of support groups on facebook-
509 the case of patients of hidradenitis suppurativa, *Multimedia Tools and*
510 *Applications* 78 (2019) 3321–3339.
- 511 [4] Y. Bengio, A. Courville, P. Vincent, Representation learning: A re-
512 view and new perspectives, *IEEE transactions on pattern analysis and*
513 *machine intelligence* 35 (2013) 1798–1828.

- 514 [5] Y. Bengio, R. Ducharme, P. Vincent, C. Jauvin, A neural probabilistic
515 language model, *Journal of machine learning research* 3 (2003) 1137–
516 1155.
- 517 [6] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed
518 representations of words and phrases and their compositionality, in:
519 *Advances in neural information processing systems*, 2013, pp. 3111–
520 3119.
- 521 [7] A.-L. Barabási, M. Pósfai, *Network science*, Cambridge University Press,
522 2016. URL: <http://barabasi.com/networksciencebook/>.
- 523 [8] Q. Wang, Z. Mao, B. Wang, L. Guo, Knowledge graph embedding: A
524 survey of approaches and applications, *IEEE Transactions on Knowledge
525 and Data Engineering* 29 (2017) 2724–2743.
- 526 [9] W. L. Hamilton, R. Ying, J. Leskovec, Representation learning on
527 graphs: Methods and applications, *arXiv preprint arXiv:1709.05584*
528 (2017).
- 529 [10] P. Cui, X. Wang, J. Pei, W. Zhu, A survey on network embedding, *IEEE
530 Transactions on Knowledge and Data Engineering* 31 (2018) 833–852.
- 531 [11] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph embedding by
532 translating on hyperplanes., in: *Aaai*, volume 14, 2014, pp. 1112–1119.
- 533 [12] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, O. Yakhnenko,
534 Translating embeddings for modeling multi-relational data, in: *Ad-
535 vances in neural information processing systems*, 2013, pp. 2787–2795.
- 536 [13] A. Bordes, X. Glorot, J. Weston, Y. Bengio, A semantic matching energy
537 function for learning with multi-relational data, *Machine Learning* 94
538 (2014) 233–259.
- 539 [14] M. Nickel, V. Tresp, H.-P. Kriegel, Factorizing yago: scalable machine
540 learning for linked data, in: *Proceedings of the 21st international con-
541 ference on World Wide Web*, 2012, pp. 271–280.
- 542 [15] M. Nickel, V. Tresp, H.-P. Kriegel, A three-way model for collective
543 learning on multi-relational data., in: *Icml*, volume 11, 2011, pp. 809–
544 816.

- 545 [16] Y. Tay, A. Luu, S. C. Hui, Non-parametric estimation of multiple embed-
546 dings for link prediction on dynamic knowledge graphs, in: Proceedings
547 of the AAAI Conference on Artificial Intelligence, volume 31, 2017.
- 548 [17] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word
549 representations in vector space, arXiv preprint arXiv:1301.3781 (2013).
- 550 [18] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: Online learning of social
551 representations, in: Proceedings of the 20th ACM SIGKDD interna-
552 tional conference on Knowledge discovery and data mining, ACM, 2014,
553 pp. 701–710.
- 554 [19] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, Q. Mei, Line: Large-scale
555 information network embedding, in: Proceedings of the 24th interna-
556 tional conference on world wide web, International World Wide Web
557 Conferences Steering Committee, 2015, pp. 1067–1077.
- 558 [20] A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks,
559 in: Proceedings of the 22nd ACM SIGKDD international conference on
560 Knowledge discovery and data mining, ACM, 2016, pp. 855–864.
- 561 [21] G. H. Nguyen, J. B. Lee, R. A. Rossi, N. K. Ahmed, E. Koh, S. Kim,
562 Continuous-time dynamic network embeddings, in: Companion Pro-
563 ceedings of the The Web Conference 2018, International World Wide
564 Web Conferences Steering Committee, 2018, pp. 969–976.
- 565 [22] P. Goyal, N. Kamra, X. He, Y. Liu, Dyngem: Deep embedding method
566 for dynamic graphs, arXiv preprint arXiv:1805.11273 (2018).
- 567 [23] L. Du, Y. Wang, G. Song, Z. Lu, J. Wang, Dynamic network embedding:
568 An extended approach for skip-gram based network embedding., in:
569 IJCAI, 2018, pp. 2086–2092.
- 570 [24] S. Mahdavi, S. Khoshraftar, A. An, dynnode2vec: Scalable dynamic
571 network embedding, in: 2018 IEEE International Conference on Big
572 Data (Big Data), IEEE, 2018, pp. 3762–3765.
- 573 [25] P. Goyal, S. R. Chhetri, A. Canedo, dyngraph2vec: Capturing network
574 dynamics using dynamic graph representation learning, Knowledge-
575 Based Systems (2019) 104816.

- 576 [26] Y. Han, J. Tang, Q. Chen, Network embedding under partial monitoring
577 for evolving networks., in: IJCAI, 2019, pp. 2463–2469.
- 578 [27] U. Singer, I. Guy, K. Radinsky, Node embedding over tem-
579 poral graphs, in: Proceedings of the Twenty-Eighth Interna-
580 tional Joint Conference on Artificial Intelligence, IJCAI-19, Interna-
581 tional Joint Conferences on Artificial Intelligence Organization, 2019,
582 pp. 4605–4612. URL: <https://doi.org/10.24963/ijcai.2019/640>.
583 doi:10.24963/ijcai.2019/640.
- 584 [28] J. R. Hurley, R. B. Cattell, The procrustes program: Producing direct
585 rotation to test a hypothesized factor structure, Behavioral science 7
586 (1962) 258–262.
- 587 [29] M. C. Gonzalez, C. A. Hidalgo, A.-L. Barabasi, Understanding individ-
588 ual human mobility patterns, nature 453 (2008) 779.
- 589 [30] R. Albert, Scale-free networks in cell biology, Journal of cell science 118
590 (2005) 4947–4957.
- 591 [31] V. Boginski, S. Butenko, P. M. Pardalos, Statistical analysis of financial
592 networks, Computational statistics & data analysis 48 (2005) 431–443.
- 593 [32] A.-L. Barabási, E. Bonabeau, Scale-free networks, Scientific american
594 288 (2003) 60–69.
- 595 [33] R. Dunford, Q. Su, E. Tamang, The pareto principle (2014).
- 596 [34] I. DRYDEN, Kv: Statistical shape analysis, John Willey, New York
597 (1999).
- 598 [35] L. Tang, H. Liu, Relational learning via latent social dimensions, in:
599 Proceedings of the 15th ACM SIGKDD international conference on
600 Knowledge discovery and data mining, ACM, 2009, pp. 817–826.
- 601 [36] L. Tang, H. Liu, Scalable learning of collective behavior based on sparse
602 social dimensions, in: Proceedings of the 18th ACM conference on
603 Information and knowledge management, ACM, 2009, pp. 1107–1116.
- 604 [37] A. Sankar, Y. Wu, L. Gou, W. Zhang, H. Yang, Dysat: Deep neural
605 representation learning on dynamic graphs via self-attention networks,

- 606 in: Proceedings of the 13th International Conference on Web Search and
607 Data Mining, 2020, pp. 519–527.
- 608 [38] T. N. Kipf, M. Welling, Semi-supervised classification with graph con-
609 volutional networks, arXiv preprint arXiv:1609.02907 (2016).
- 610 [39] S. Nandanwar, M. N. Murty, Structural neighborhood based classifica-
611 tion of nodes in a network, in: Proceedings of the 22nd ACM SIGKDD
612 International Conference on Knowledge Discovery and Data Mining,
613 ACM, 2016, pp. 1085–1094.
- 614 [40] A. K. McCallum, K. Nigam, J. Rennie, K. Seymore, Automating the
615 construction of internet portals with machine learning, *Information Re-
616 trieval* 3 (2000) 127–163.