

Structure and Features Fusion with Evidential Graph Convolutional Neural Network for Node Classification

Anonymous ACL submission

Abstract

1 Recently, text-enhanced network
2 representation learning has achieved great
3 success by taking advantage of rich text
4 information and network structure
5 information. However, content-rich
6 network representation learning and
7 quantifying classification uncertainty are
8 challenging when it comes to integrating
9 complex structural dependencies and rich
10 content features at an evidence level. In this
11 paper, we propose an evidential graph
12 representation learning model (EGCN),
13 which can not only fuse network structure
14 and content information into a more
15 complete and powerful representation for
16 each node, but also assess the quality of
17 graph node features to improve
18 classification accuracy. To achieve better
19 fusion, we integrate the node's features
20 representation into structure-aware
21 representation through a delivery operator.
22 Besides, to overcome the difficulty of
23 predicting node classification confidence,
24 we employ a novel module based on
25 Dirichlet distribution theory of evidence
26 and subject opinion learning to collect the
27 evidence of the class probabilities.
28 Experimental results on three real-world
29 networks show that our model can improve
30 both node classification accuracy and
31 robustness as compared to all baselines.
32

1 Introduction

34 Content-rich networks are graphs with node
35 features and network structures widely applied in
36 academic citation networks, recommendation
37 systems, etc. However, because of the complex
38 non-Euclidean graph structure, capturing structure
39 and feature information is a challenging task on
40 machine learning approaches.

41 Recently, research on analyzing networks with
42 deep learning has received widespread attention. In
43 particular, graph convolutional networks (GCNs)
44 (Defferrard et al., 2016; Kipf and Welling 2017),
45 which obtain node embeddings through the
46 propagation and aggregation of the features on
47 network topology, have achieved great success.
48 While the success of GCNs and their variants, a key
49 issue with them is that the accuracy of multi-typed
50 features classification varies greatly. To address
51 this problem, (Yang et al., 2015) proposed text-
52 associated DeepWalk (TADW), which
53 incorporates text features of vertices into network
54 representation learning under the framework of
55 matrix factorization. But the model can only handle
56 the text attributes. (Cui et al., 2020) presented an
57 adaptive graph encoder (AGE), a novel attributed
58 graph embedding framework which applied a
59 carefully-designed Laplacian smoothing filter.
60 Nevertheless, all the above methods are designed
61 to handle the single-typed feature and the node
62 features only serve as an initial solution of
63 embeddings.

64 GCNs obtain node features from a local mixing
65 state of propagation by limiting the number of
66 propagations to two or three layers. However, this
67 will further make GCNs rely heavily on the local
68 homophily of topology, that is, neighborhoods
69 should be similar, a very strong assumption in
70 many real-world text-rich networks. A lot of
71 methods are developed to handle the topological
72 limitations of GCNs. For example, several studies
73 attempted to utilize the self-supervised learning
74 methods (Zhu et al., 2020; You et al., 2020) which
75 use several highly credible labels derived from
76 GCNs to optimize the topological channels in the
77 following propagation of GCNs. These existing
78 methods have achieved reasonable results at
79 handling the topological limitations of GCNs and
80 thus improved the performance of GCNs. However,
81 from the level of model architectures, an ideal way

82 may be that the convolutions of features (on
83 topology) and topology (on features) play together
84 in the same system. By jointly training the BERT
85 and GCN modules within Bert-GCN (Lin et al.,
86 2021), this model is able to leverage the advantages
87 of both worlds: large-scale pretraining which takes
88 the advantage of the massive amount of raw data
89 and transductive learning. However, the model was
90 trained with the BERT feature of node text and it
91 cannot utilize multi-typed features like term fre-
92 quency-inverse document frequency (TF-IDF) or
93 SimCSE(Gao et al., 2021).

94 To the best of our knowledge, no work has been
95 devoted to exploring both the multi-typed features
96 and the semantic graph relationships in an efficient
97 way. In this paper, we develop a unified deep model
98 (EGCN) to capture both text-rich information and
99 topology structure features. The training process of
100 EGCN consists of three parts which focus on
101 preserving the information of multi-typed features,
102 network topologies, and node classification
103 confidence, respectively. The graph structure
104 information is mined by modeling the first-order
105 and second-order proximities between nodes; the
106 text features are processed with TF-IDF, BERT,
107 and SimCSE methods to capture the different
108 patters. Also, our model combines different node
109 features at an evidence level, which produces a
110 stable and reasonable uncertainty estimation.
111 Figure 1 shows the illustration of our
112 implementation of EGCN. The main contributions
113 of this work are summarized as follows.

114 (1) We propose a unified deep model (EGCN) to
115 learn the embedding vector for each node of the
116 network by considering both multi-typed features
117 and the graph semantic relationships,
118 simultaneously.

119 (2) We develop a novel multi-typed features
120 classification method aiming to provide trusted and
121 interpretable decisions in an effective and efficient
122 way.

123 (3) We run extensive experiments which validate
124 the superior accuracy and robustness of our model
125 thanks to the promising uncertainty estimation and
126 multi-typed features integration strategy.

127 2 Related Work

128 2.1 Enriching Graph Embeddings with 129 External Text

130 Graph convolutional networks (GCN) is
131 connectionist models that fusion dependencies and

132 relations between graph nodes (Hamilton et al.,
133 2017; Keyulu et al., 2018). Besides structural
134 properties, nodes in a graph are often affiliated with
135 various contents, such as abstract or title text in the
136 academic citation network. Such networks are
137 called text-rich networks, and have been
138 extensively studied (Li et al., 2017; Zhang et al.,
139 2018; Zhou et al., 2018; Velićkovic' et al. 2019;
140 Meng et al. 2019). Their goal is to preserve not only
141 the network structure, but also the node attribute
142 proximity in learning representations. Recently,
143 much efforts have been made to gain insights from
144 attributed networks (Liao et al., 2018; Yang et al.,
145 2018; Li et al., 2017). Some approaches (Huang et
146 al., 2017; Xiao et al., 2017) simply take the label
147 information into consideration, while others utilize
148 more detailed attribute information. The key point
149 of attributed network embedding lies in simultane-
150 ously capturing node attributes, network
151 structure and their relation-ship into hidden
152 representations. Our work is inspired by the work
153 of using graph neural networks to fusion node
154 features (Zhang et al., 2020). Existing works that
155 combine BERT and GNNs uses graph to model
156 relationships between tokens within a single
157 document sample (Lu et al., 2020), which fall into
158 the category of inductive learning. But different
159 from these works, we focus on combining, and
160 show that multi-typed features can significantly
161 benefit from uncertainty-based learning model.

162 2.2 Uncertainty-based Learning

163 The history of learning uncertainty-aware
164 predictors is concurrent with the advent of modern
165 Bayesian approaches to machine learning.
166 Bayesian neural networks (BNNs) (Neal et al.,
167 2012) endow deep models with uncertainty by
168 replacing the deterministic weight parameters with
169 distributions. Because BNNs need to consume a lot
170 of computing power when performing inference
171 calculations, a more stable and effective method,
172 MC-dropout (Gal et al., 2016), was proposed. The
173 inference calculation in this model is done by
174 dropout sampling from the training and test
175 weights. Ensemble based methods (Lak-
176 shminarayanan et al., 2017) train and integrate
177 multiple deep networks and also achieve promising
178 performance. Instead of indirectly modeling
179 uncertainty through net-work weights, the
180 algorithm (Sensoy et al., 2018) introduces the
181 subjective logic theory to directly model
182 uncertainty without ensemble or Monte Carlo

183 sampling. Building upon RBF networks, the
 184 distance between test samples and prototypes can
 185 be used as the agency for deterministic uncertainty
 186 (van Amersfoort et al., 2020). Benefiting from the
 187 learned weights of different tasks with
 188 homoscedastic un-certainty learning, (Kendall et
 189 al., 2018) achieves impressive performance in
 190 multi-task learning. (Han et al. 2021) utilized
 191 multiple views to promote both classification reli-
 192 ability and robustness by integrating evidence from
 193 each view. Dempster-Shafer Evidence Theory is
 194 about the theoretical method of the confidence
 195 function, which directly models uncertainty. DST
 196 allows beliefs from different sources to be
 197 combined with various fusion operators to obtain a
 198 new belief that considers all available evidence
 199 (Jøsang et al., 2012).

200 3 EGCN Model

201 3.1 Problem Statement

202 A text-rich network is a network $G = \{V, R, X\}$,
 203 where R is the set of relations. V is the set of nodes.
 204 X is a matrix that encodes node attributes
 205 information for n nodes. Given an attributed
 206 network G and the set of adjacency matrices A , the
 207 task of structure and features fusion with EGCN for
 208 node classification is to learn the multi-type
 209 features, and determine the confidence of each type
 210 of feature for the node classification.

211 3.2 GCN Module

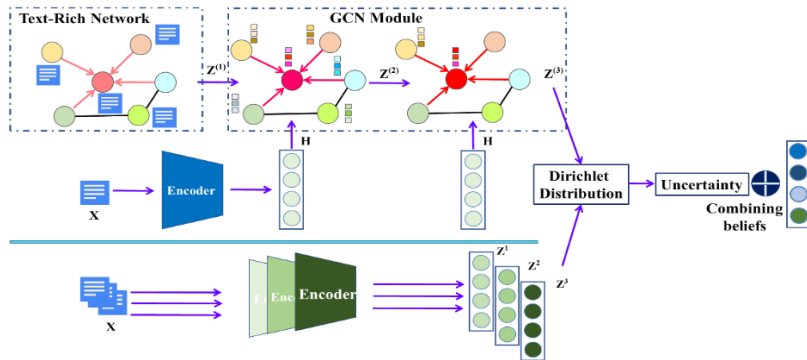
212 The pre-training models can obtain the expression
 213 of the semantic information of the text. In the
 214 section, we will introduce how to use the GCN
 215 module to propagate these representations
 216 generated by the pre-training models. Once all the
 217 multi-typed representations learned by BERT and
 218 SimCSE are integrated into GCN, then the output
 219 feature of GCN will be able to accommodate for
 220 two different kinds of information, i.e., data itself
 221 and relationship. In particular, with the weight
 222 matrix W , the representation learned by the ℓ -th
 223 layer of GCN, $Z^{(\ell)}$, can be obtained by the
 224 following convolutional operation:

$$225 \quad Z^{(\ell)} = \phi \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} Z^{(\ell-1)} W^{(\ell-1)} \right) \quad (1)$$

226 where $\tilde{A} = A + I$, $\phi(\cdot)$ is an activation function such
 227 as ReLU. As can be seen from Eq. 1, the
 228 representation $Z^{(\ell-1)}$ will propagate through the
 229 normalized adjacency matrix $\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$ to obtain
 230 the new representation $Z^{(\ell)}$. Considering that the
 231 representation learned by BERT or SimCSE, we
 232 combine the two representations $Z^{(\ell-1)}$ and H
 233 together to get a more complete and powerful
 234 representation as follows:

$$235 \quad \tilde{Z}^{(\ell-1)} = (1 - \epsilon)Z^{(\ell-1)} + \epsilon H \quad (2)$$

236 where ϵ is a balance coefficient, H represents text
 237 feature.



238

239 Figure 1 Visual illustration of our implementation of EGCN. Text-rich network obtains the graph structure
 240 information through GCN module, and at the same time, the text features of encoders such as BERT or SimCSE
 241 are fused with GCN with Eq. 2. The text features such as TF-IDF, BERT and SimCSE and the features output by
 242 the GCN module are mapped using Dirichlet distribution to calculate the confidence and uncertainty of each
 243 feature through Uncertainty and the Theory of Evidence, combining them with Dempster-Shafer theory, and
 244 obtain the confidence and classification uncertainty after all the features are combined.

245 3.3 Uncertainty and the Theory of Evidence

246 The Dempster-Shafer Theory of Evidence (DST)
 247 is a theory on belief functions, was first proposed
 248 by Dempster and is a generalization of the

249 Bayesian theory to subjective probabilities.
 250 Subjective Logic (SL) formalizes DST's notion of
 251 belief assignments over a frame of discernment as
 252 a Dirichlet Distribution. Hence, we introduce a
 253 principle of evidential theory-based uncertainty

estimation technique which can provide more accurate uncertainty and allow us to flexibly integrate multi-typed features for trusted classification decision making. More specifically, SL considers a frame of K mutually exclusive singletons (e.g., class labels) by providing a belief mass b_k for each singleton $k = 1, \dots, K$ and providing an overall uncertainty mass of u . These $K + 1$ mass values are all non-negative and sum up to one. Note that evidence e_k refers to the metrics collected from the multi-typed features to support the classification in Figure 2.

$$u + \sum_{k=1}^K b_k = 1 \quad (3)$$

where $u \geq 0$ and $b_k \geq 0$ for $k = 1, \dots, K$.

A belief mass b_k for a singleton k is computed using the evidence for the singleton. Let $e_k \geq 0$ be the evidence derived for the k th singleton, then the belief b_k and the uncertainty u are computed as:

$$b_k = \frac{e_k}{S} \text{ and } u = \frac{K}{S} \quad (4)$$

where $S = \sum_{k=1}^K (e_k + 1)$.

Eq. 4 actually describes the phenomenon where the more evidence observed for the k -th category, the greater the probability assigned to the k -th class. A belief mass assignment, i.e., subjective opinion, corresponds to a Dirichlet distribution with parameters $\alpha_k = e_k + 1$. That is, a subjective opinion can be derived easily from the parameters of the corresponding Dirichlet distribution using $b_k = (\alpha_k - 1)/S$. However, a Dirichlet distribution parametrized over evidence represents the density of each such probability assignment; hence it models second-order probabilities and uncertainty (Han et al. 2021). The Dirichlet distribution is a probability density function for possible values of the probability mass function p . It is characterized by K parameters $\alpha = [\alpha_1, \dots, \alpha_K]$ and is given by

$$D(p | \alpha) = \begin{cases} \frac{1}{B(\alpha)} \prod_{i=1}^K p_i^{\alpha_i-1} & \text{for } p \in S_K \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where S_K is the K -dimensional unit simplex, $B(\alpha)$ is the K -dimensional multinomial beta function.

$$S_K = \{p | \sum_{i=1}^K p_i = 1 \text{ and } 0 \leq p_1, \dots, p_K \leq 1\} \quad (6)$$

$$\mathcal{L}_{ace}(\alpha_i) = \int [\sum_{j=1}^K -y_{ij} \log(p_{ij})] \frac{1}{B(\alpha_i)} \prod_{j=1}^K p_{ij}^{\alpha_{ij}-1} dp_i = \sum_{j=1}^K y_{ij} (\psi(S_i) - \psi(\alpha_{ij})) \quad (9)$$

where $\psi(\cdot)$ is the digamma function.

Eq. 9 is the integral of the cross-entropy loss function on the simplex determined by α_i . The above loss function ensures that the correct label of

In this paper, we design and train neural networks to form multi-view opinions for the classification of a given sample i as a Dirichlet distribution $D(p_i | \alpha_i)$, where p_i is a simple representing class assignment probabilities.

3.4 Dempster's Rule of Combination for Multi-typed features Classification

Having introduced evidence and uncertainty for the single-feature case, we use the Dempster-Shafer theory of evidence to combine multi-typed features arriving at a degree of belief (represented by a mathematical object called the belief function) that focus on all the available evidence (see Figure 2). Specifically, we combine V independent sets of probability mass assignments $\{M^v\}_1^V$, where $M^v = \{b_k^v\}_{k=1}^K, u^v\}$, b refers to the confidence probability, v is the feature type, k is the node category, and u^v is the uncertainty of the node classification for the feature v . The combined calculation of confidence and uncertainty among multiple types of features is as follows:

$$\mathcal{M} = \mathcal{M}^1 \oplus \mathcal{M}^2 \quad (7)$$

The more specific calculation rule can be formulated as follows:

$$b_k = \frac{1}{1-C} (b_k^1 b_k^2 + b_k^1 u^2 + b_k^2 u^1), u = \frac{1}{1-C} u^1 u^2 \quad (8)$$

where $C = \sum_{i \neq j} b_i^1 b_j^2$ is a measure of the amount of conflict between the two mass sets, and the scale factor $1-C$ is used for normalization.

3.5 Learning to Form Opinions

The loss over a batch of training samples can be computed by summing the loss for each sample in the batch. During training, the model mine patterns in the node text and generate evidence for specific class labels based on these patterns to minimize the overall loss. For our model, given the evidence of the i -th sample obtained through the evidence network, we can get the parameter α_i (i.e., $\alpha_i^v = e^i + 1$) of the Dirichlet distribution and form the multinomial opinions $D(p_i | \alpha_i)$. we can treat $D(p_i | \alpha_i)$ as a prior on the likelihood $\text{Mult}(y_i | p_i)$ and obtain the negated logarithm of the marginal likelihood by integrating out the class probabilities

each sample generates more evidence than other classes, however, it cannot guarantee that less evidence will be generated for incorrect labels. That is to say, in our model, we expect the evidence

345 for incorrect labels to shrink to 0. To this end, the
 346 following KL divergence term is introduced:

$$\begin{aligned}
 & KL[D(p_i | \tilde{\alpha}_i) \parallel D(p_i | 1)] = & (10) \\
 & \log \left(\frac{\Gamma(\sum_{k=1}^K \tilde{\alpha}_{ik})}{\Gamma(K) \prod_{k=1}^K \Gamma(\tilde{\alpha}_{ik})} \right) + \\
 & \sum_{k=1}^K (\tilde{\alpha}_{ik} - 1) \left[\psi(\tilde{\alpha}_{ik}) - \psi \left(\sum_{j=1}^K \tilde{\alpha}_{ij} \right) \right]
 \end{aligned}$$

347 where $\tilde{\alpha}_i = y_i + (1 - y_i) \odot \alpha_i$ is the adjusted parameter
 348 of the Dirichlet distribution which can avoid
 349 penalizing the evidence of the ground-truth class
 350 to 0, and $\Gamma(\cdot)$ is the gamma function.

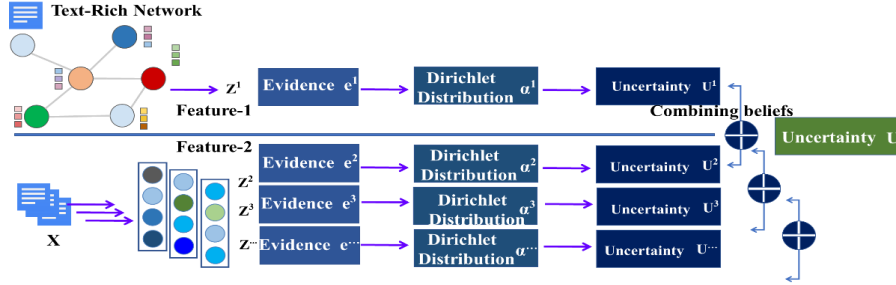
351 Let us consider that a Dirichlet distribution with
 352 zero total evidence, i.e., $S = K$, corresponds to the
 353 uniform distribution and indicates total
 354 uncertainty, i.e., $u = 1$. We achieve this by
 355 incorporating a Kullback-Leibler(KL) divergence

356 term into our loss function that regularizes our
 357 predictive distribution by penalizing those
 358 divergences from the "I do not know" state that do
 359 not contribute to data fit. The loss with this
 360 regularizing term reads:

$$361 \mathcal{L}(\alpha_i) = \mathcal{L}_{ace}(\alpha_i) + \lambda_t KL[D(p_i | \tilde{\alpha}_i) \parallel D(p_i | 1)] \quad (11)$$

362 where $\lambda_t = \min(1.0, t/10) \in [0, 1]$ is the annealing
 363 coefficient, t is the index of the current training
 364 epoch, $D(p_i | h_1, \dots, h_i)$ is the uniform Dirichlet
 365 distribution, to prevent the network from paying
 366 too much attention to the KL divergence in the
 367 initial stage of training. To ensure that all views
 368 can simultaneously form reasonable opinions and
 369 thus improve the overall opinion, we use a multi-
 370 task strategy with following overall loss function:

$$371 \mathcal{L}_{overall} = \sum_{i=1}^N [\mathcal{L}(\alpha_i) + \sum_{v=1}^V \mathcal{L}(\alpha_i^v)] \quad (12)$$



372

373 Figure 2 Illustration of trusted multi-typed features classification. The evidence of each feature is obtained using
 374 BERT, TF-IDF, SimCSE and GCN in Figure 1. The obtained evidence parameterizes the Dirichlet distribution to
 375 induce the classification probability and uncertainty. The overall uncertainty and classification probability are
 376 inferred by combining the beliefs of multiple views based on the DST. The combination rule and an example are
 377 shown in Eq.7 and Eq. 8, respectively.

378 4 Experiments

379 In this section, we run experiments on three real-
 380 world da-tasets: Cora, Citeseer and DBLP. We
 381 compare EGCN to the following models: BERT
 382 (Devlin et al. 2018), SimCSE (Gao et al. 2021),
 383 GCN (Kipf et al. 2017), GAT (Veličković et al.
 384 2017), GraphSage (Hamilton et al. 2017), TADW
 385 (Yang et al. 2015), BertGCN (Lin et al. 2021) to
 386 demonstrate the effectiveness of proposed model.
 387 We also prove EGCN model can produce trusted
 388 classification decisions on different types of
 389 attributed information.

390 4.1 Datasets

391 Cora data is an open citation network data set,
 392 containing 7 types of papers. The network contains
 393 2211 paper nodes and 5214 citation relationships.
 394 Each paper contains an average of 169 words, and
 395 the vocabulary of the entire data set contains a total

396 of 12619 words. The Citeseer data set consists of
 397 papers from 10 interdisciplinary research fields, it
 398 contains 4610 nodes and 5923 edges. The DBLP
 399 data set is a comprehensive data set covering 4
 400 types of papers, the network contains 13,404 nodes
 401 and 39861 edges.

Dataset	Cora	Citeseer	DBLP
# Nodes	2211	4610	13404
# Edge	5214	5923	39861
# Text	169	10	10
# Classes	7	10	4

402

Table 1 Dataset statistics

403 Table 1 illustrates the details of datasets used in
 404 our experiment. #Text denotes the average number
 405 of words contained in each text node

406 4.2 Experiment Setups

407 For all methods using the BERT model, we use
 408 BERT-base architecture with pre-trained weights
 409 from the original authors and adapted by

410 HuggingFace Transformers library³. We then fine-
 411 tune it using masked language model objective on
 412 the three real-world datasets: Cora, Citeseer and
 413 DBLP with a 10^{-5} learning rate. We set the number
 414 of layers to 2, And the hidden layer dimension is
 415 equal to 768, in order to be consistent with the
 416 dimension of the graph structure data and the
 417 BERT feature data. For the SimCSE method in the
 418 article, the temperature constant of the contrast loss
 419 function is set to 0.05.

420 As for our model, for different data sets, because
 421 the number of words in the document is different,
 422 the ability to represent different text features is not
 423 the same. Therefore, it is necessary to effectively
 424 fuse different features with confidence. The main
 425 idea is to choose the single-typed feature with
 426 higher classification capabilities. We choose multi-
 427 typed features from BERT, TF-IDF, and SimCSE
 428 methods and incorporate GCN output features to
 429 obtain the confidence and uncertainty of each
 430 feature for classification. In the fusion processing
 431 of multi-typed features, we set the fusion ratio λ to
 432 be 0.3.

433 4.3 Main Results

434 Table 2 presents the test accuracy of each model.
 435 We can see that EGCN outperforms other models
 436 in the three data sets. For the pre-training features
 437 of text data, the accuracy on the three data sets is
 438 low, the accuracy on the GCN, GAT and
 439 GraphSage models is improved to a certain extent.
 440 Our method has a higher advantage. Competing
 441 with the strongest baseline BertGCN, our model
 442 outperforms it by 3% on Cora, by 9% on Citeseer,
 443 by 6% on DBLP.

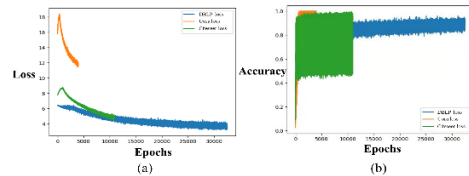
Models	Cora	Citeseer	DBLP
BERT	0.55	0.60	0.63
SimCSE	0.71	0.67	0.67
GCN	0.78	0.72	0.65
GAT	0.79	0.80	0.69
GraphSage	0.78	0.81	0.71
BertGCN	0.83	0.88	0.84
EGCN	0.86	0.97	0.90

444 Table 2: Experimental results of node classification.

445 From the text, graph structure, text and graph
 446 structure fusion of Table 2 to multi-typed features
 447 trusted node classification, the following
 448 conclusions can be drawn: for graph node
 449 classification, Algorithms using both feature and
 450 graph information achieve better performance than
 451 methods leveraging information from single source.
 452 This investigation demonstrates that features and

453 graph structure contribute to classification from
 454 different perspectives.

455 Figure 3 shows the loss (a) and the corresponding
 456 accuracy (b) of the model training process on the three
 457 data sets. It can be seen from the figure that due to the
 458 large amount of DBLP data, the loss and accuracy curve
 459 is longer. At the same time, (b) reflects the convergence
 460 of our model and shows higher accuracy. For the
 461 accuracy curve of (b), there is a jitter phenomenon,
 462 which can be explained by (d) and (h) in Figure 4. The
 463 classification error is caused by the conflict of multiple
 464 types of features.



466 Figure 3 Illustration of the training process.

467 4.4 Ablation study

468 Table 3 shows the classification results of text
 469 features. For the three data sets, learn feature
 470 expressions of the node's TF-IDF, SimCSE, and
 471 BERT. BERT768 refers to the out-put dimension is
 472 768, and BERT512 and BERT256 refer to the slave
 473 the first 512-dimensional and 256-dimensional
 474 features cropped from the 768-dimensional
 475 features. From Table 3, it can be seen that for node
 476 text classification, the TF-IDF feature performs the
 477 best for the classification results of the three data
 478 sets.

MLP	Cora	Citeseer	DBLP
TF-IDF	0.85	0.84	0.79
SimCSE	0.71	0.67	0.67
Bert768	0.55	0.6	0.63
Bert512	0.49	0.41	0.58
Bert256	0.29	0.33	0.5

479 Table 3 Experimental results of node text multi-type
 480 features classification.

GCN+MLP	Cora	Citeseer	DBLP
TF-IDF	0.8381	0.9397	0.8715
SimCSE	0.6524	0.7397	0.6044
Bert768	0.7841	0.7511	0.5288
Bert512	0.7984	0.8519	0.8265
Bert256	0.7714	0.8153	0.7892

481 Table 4 Experimental results of classification of GCN
 482 nodes with multiple types of features

483 Table 4 shows the node classification results.
 484 We utilize the 2-layer GCN to aggregate the node
 485 feature neighbor in-formation after the text node
 486 has been learned by different feature expressions.
 487 It can be seen that, except for the feature vector

expressed by SimCSE, the node classification accuracy obtained a big improvement.

Table 5 shows that different features are linearly fused when GCN is used for feature aggregation. $\lambda=1$ means the node feature classification using GCN, and $\lambda=0$ means the node text feature classification is used. As can be seen from the table, for the three data Collection, linear fusion between multiple features, the output features are classified, and the accuracy is greatly improved.

λ GCN+(1- λ)	Cora	Citeseer	DBLP
Features			
TF-IDF	0.8048	0.8785	0.8436
SimCSE	0.7952	0.8715	0.8407
Bert768	0.8126	0.8708	0.8467
Bert512	0.7952	0.8769	0.8441
Bert256	0.7984	0.8769	0.8449

Table 5 Experimental results of node classification of multi-type features with linear fusion of text features and structural features

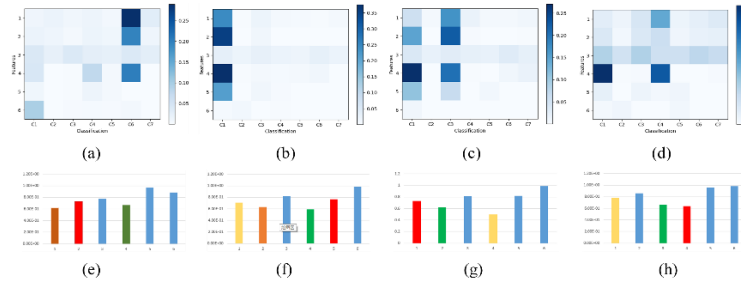


Figure 4 Illustration of node multi-type features classification confidence

501
502

For the progressive experimental results in Tables 3, 4, 5 and Table 2, we have the following observations:

(1) The basic observation is that our proposed EGCN framework achieves better results on three datasets compared with baseline methods and variants. This shows the effectiveness of our proposed model in modeling node features and network topology. By comparing EGCN with baseline methods, we can further infer the advantage of aggregating the multi-type features of nodes and structures to classify the confidence of nodes.

(2) From node text feature classification, text feature input to GCN node classification, text feature and GCN structure feature linear aggregation node feature classification to text feature, text feature and graph structure feature multi-type features for category confidence classification, classification The accuracy continues to improve, showing the effective fusion of multiple types of features and the experimental verification of the confidence of each type of feature on the classification from different types of features.

In Figure 4, for the four samples of the Cora dataset, the rows in Figures (a) (b) (c) (d) represent six different types of features (including the text features and structural features of the graph nodes), and the list is up to 7 classification types, color intensity codes, the confidence of each feature for the 7 types of samples, the darker the color, the more confidence, the sixth column of Figure (a) has a darker color, expressing the confidence of multiple types of features for the sixth category, Figure (b) is a darker color in the first column. Express the confidence of multiple types of features for the first type. The colors in the first and third columns in Figure (c) are darker, but the color depth of the third type is evenly distributed. The experimental results show that it is classified as the

third type. The color distribution in Figure (d) is messy and the final classification is wrong. You can see the third row of main features are all dark in color, expressing that this type of feature is confident in the seven classification results, and the dark colors in the first, third, and fourth columns cause confidence confusion, and the experimental results are wrong. Figure (e) (f) (g) (h) represents the classification uncertainty of 6 different types of features for 4 samples. The lower the histogram, the lower the uncertainty and the higher the certainty. Figure (e) (f) (g) (h) respectively correspond to Figures (a) (b) (c) (d), in which the uncertainty of the six features for sample classification can also effectively reflect the confidence of the classification results.

5 Conclusion

In this work, we propose a novel trusted multi-typed features graph node classification (EGCN) model which, based on the Dempster-Shafer evidence theory, can produce trusted classification decisions on multi-typed features and can jointly learn low-dimensional representations of both nodes and features for text-rich networks. Our algorithm focuses on decision-making by fusing the uncertainty of multi-typed features, which is essential for making trusted decisions. Furthermore, our model can produce the uncertainty of a current decision while making the final classification, providing interpretability. The empirical results validate the effectiveness of the proposed algorithm in classification accuracy.

570 References

- 571 Kipf, T. N.; Welling, M. 2017. Semi-Supervised
572 Classification with Graph Convolutional Networks.
573 In ICLR.
- 574 Defferrard, M; Bresson, X.; Vandergheynst, P. 2016.
575 Convolutional neural networks on graphs with fast
576 localized spectral filtering. *Advances in neural
577 information processing systems*, 3844-3852.
- 578 Yang, C.; Liu, Z., Zhao, D., Sun, M.; Chang, E. 2015.
579 Network representation learning with rich text
580 information. In *Twenty-fourth international joint
581 conference on artificial intelligence*.
- 582 Cui, G.; Zhou, J.; Yang, C.; Liu, Z. 2020. Adaptive
583 graph encoder for attributed graph embedding. In
584 *Proceedings of the 26th ACM SIGKDD
585 International Conference on Knowledge Discovery
586 & Data Mining*, 976-985.
- 587 Zhu, Q.; Du, B.; Yan, P. 2020. Self-supervised training
588 of graph convolutional networks. *arXiv preprint
589 arXiv:2006.02380*.
- 590 You, Y.; Chen, T.; Wang, Z.; Shen, Y. 2020. When does
591 self-supervision help graph convolutional
592 networks?. In *International Conference on Machine
593 Learning* (pp. 10871-10880). PMLR.
- 594 Lin Y; Meng Y; Sun X. 2021. BertGCN: Transductive
595 Text Classification by Combining GNN and
596 BERT. In *Findings of the Association for
597 Computational Linguistics: ACL-IJCNLP*, 1456-
598 1462.
- 599 Gao, T.; Yao, X.; Chen, D. 2021. SimCSE: Simple
600 Contrastive Learning of Sentence Embeddings. In
601 *Empirical Methods in Natural Language Processing*.
602 In EMNLP.
- 603 Hamilton W; Zhitao Y; and Jure L. 2017. Inductive
604 representation learning on large graphs. In
605 *Advances in neural information processing
606 systems*, 1024-1034.
- 607 Keyulu X; Weihua H; Jure L; and Stefanie J. 2018.
608 How powerful are graph neural networks? *arXiv
609 preprint arXiv:1810.00826*.
- 610 Li, J.; Dani, H.; Hu, X.; Tang, J.; Chang, Y.; and Liu,
611 H. 2017. Attributed network embedding for learning
612 in a dynamic environment. In *CIKM*. ACM.
- 613 Zhang, Z.; Yang, H.; Bu, J.; Zhou, S.; Yu, P.; Zhang, J.;
614 Ester, M.; and Wang, C. 2018. Anrl: Attributed
615 network representation learning via deep neural
616 networks. In *IJCAI*.
- 617 Zhou, S.; Yang, H.; Wang, X.; Bu, J.; Ester, M.; Yu, P.;
618 Zhang, J.; and Wang, C. 2018. Prre: Personalized
619 relation ranking embedding for attributed networks.
620 In *CIKM*. ACM.
- 621 Veličković, P.; Fedus, W.; Hamilton, W. L.; Li, P.;
622 Bengio, Y.; and Hjelm, R. D. 2019. Deep graph
623 infomax. *ICLR*.
- 624 Meng, Z.; Liang, S.; Bao, H.; and Zhang, X. 2019. Co-
625 embedding attributed networks. In *WSDM*. ACM.
- 626 Liao L.; He X.; Zhang H. 2018. Attributed social
627 network embedding, *IEEE Transactions on
628 Knowledge and Data Engineering*.
- 629 Yang H.; Pan, S; Zhang P; Chen L.; Lian D.; and Zhang
630 C.. 2018. Binarized attributed network embedding,”
631 in *ICDM*. IEEE, 1476-1481.
- 632 Li J.; Dani H.; Hu X.; Tang J.; Chang Y.; and Liu H.
633 2017. Attributed network embedding for learning in
634 a dynamic environment. In *CIKM*. ACM, 387-396.
- 635 Huang X.; Li J.; and Hu X. 2017. Accelerated
636 attributed network embedding. In *SIAM
637 International Conference on Data Mining*, 633-641.
- 638 Xiao H.; Jundong L.; and Xia H. 2017. Label informed
639 attributed network embedding. In *WSDM*. ACM,
640 731-739.
- 641 Zhang H and Zhang X. 2020. Text graph transformer
642 for document classification. In *Proceedings of the
643 2020 Conference on Empirical Methods in Natural
644 Language Processing (EMNLP)*, 8322-8327.
- 645 Lu Z; Pan Du P; and Nie J. 2020. Vgcn-bert:
646 augmenting bert with graph embedding for text
647 classification. In *European Conference on
648 Information Retrieval*, pages 369-382. Springer.
- 649 Radford M Neal. 2012 *Bayesian learning for neural
650 networks*. Springer Science & Business Media.
- 651 Gal Y and Ghahramani Z. 2016. Dropout as a bayesian
652 approximation: Representing model uncertainty in
653 deep learning. In *International Conference on
654 Machine Learning*, 1050-1059.
- 655 Lakshminarayanan B; Pritzel A; and Blundell C.
656 2017. Simple and scalable predictive uncertainty
657 estimation using deep ensembles. In *Advances in
658 Neural Information Processing Systems*, 6402-6413.
- 659 Sensoy M; Kaplan L and Kandemir M. 2018.
660 Evidential deep learning to quantify classification
661 uncertainty. In *Advances in Neural Information
662 Processing Systems*, 3179-3189.
- 663 van Amersfoort J; Smith L; Whyte Te ,Y and Gal Y.
664 2020. Uncertainty estimation using a single deep
665 deterministic neural network. In *International
666 Conference on Machine Learning*.
- 667 Alex Kendall, Yarin Gal, and Roberto Cipolla. 2018.
668 Multi-task learning using uncertainty to weigh
669 losses for scene geometry and semantics. In
670 *Proceedings of the IEEE Conference on Computer
671 Vision and Pattern Recognition*, 7482-7491.

- 672 Han, Z., Zhang, C., Fu, H., & Zhou, J. T. 2021. Trusted
673 Multi-View Classification. In ICLR.
- 674 Jøsang A and Hankin R. 2012. Interpretation and fusion
675 of hyper opinions in subjective logic. In 2012 15th
676 International Conference on Information
677 Fusion,1225–1232
- 678 Devlin, J.; Chang, M. W.; Lee, K; Toutanova, K. 2018.
679 Bert: Pre-training of deep bidirectional transformers
680 for language under-standing. arXiv preprint
681 arXiv:1810.04805.
- 682 Veličković, P., Cucurull, G., Casanova, A., Romero, A.,
683 Lio, P., & Bengio, Y. 2017. Graph attention
684 networks. arXiv preprint arXiv:1710.10903.
- 685 Hamilton W, Ying Z, Leskovec J.2017.Inductive
686 representation learning on large graphs. Advances in
687 Neural Information Processing Systems, 1024-1034.