

Joint Inference for Aspect-Level Sentiment Analysis by Deep Neural Networks and Linguistic Hints

Yanyan Wang¹, Qun Chen¹, Murtadha Ahmed, Zhanhuai Li, Wei Pan, and Hailong Liu

Abstract—The state-of-the-art techniques for aspect-level sentiment analysis focused on feature modeling using a variety of deep neural networks (DNN). Unfortunately, their performance may still fall short of expectation in real scenarios due to the semantic complexity of natural languages. Motivated by the observation that many linguistic hints (e.g., sentiment words and shift words) are reliable polarity indicators, we propose a joint framework, SenHint, which can seamlessly integrate the output of deep neural networks and the implications of linguistic hints in a unified model based on Markov logic network (MLN). SenHint leverages the linguistic hints for multiple purposes: (1) to identify the easy instances, whose polarities can be automatically determined by the machine with high accuracy; (2) to capture the influence of sentiment words on aspect polarities; (3) to capture the implicit relations between aspect polarities. We present the required techniques for extracting linguistic hints, encoding their implications as well as the output of DNN into the unified model, and joint inference. Finally, we have empirically evaluated the performance of SenHint on both English and Chinese benchmark datasets. Our extensive experiments have shown that compared to the state-of-the-art DNN techniques, SenHint can effectively improve polarity detection accuracy by considerable margins.

Index Terms—Deep neural networks, linguistic hints, aspect-level sentiment analysis

1 INTRODUCTION

ASPECT-LEVEL sentiment analysis (ALSA) [1], a fine-grained classification task, has recently become an active research area in NLP. Its goal is to extract the opinions expressed towards different aspects of a product. ALSA can provide important insights into products to both consumers and businesses [2]. In the literature [3], two finer subtasks of ALSA have been studied: aspect-category sentiment analysis (ACSA) and aspect-term sentiment analysis (ATSA). ACSA aims to predict the sentiment polarity towards a few predefined aspect categories, which may not explicitly appear in the text. ATSA instead deals with explicit aspects involving a single word or a multi-word phrase. In this paper, we target both ACSA and ATSA. Consider the running example shown in Table 1, in which R_i and S_{ij} denote the review and sentence identifiers respectively. It can be observed that in R_2 , the aspect term “battery” explicitly appears in the sentence S_{21} , while the sentence S_{22} does not explicitly contain its target aspect term (“laptop#performance”). ACSA has to detect the polarities of the aspects in both S_{21} and S_{22} . In contrast, ATSA only needs to detect the aspect polarity in S_{21} .

The state-of-the-art solutions for aspect-level sentiment analysis [4], [5] are mainly built on a variety of deep neural networks (DNN), which can automatically learn multiple levels of feature representation. Even though the DNN techniques can achieve empirically better performance than the previous alternatives (e.g., the techniques based on lexicon [6], [7] and SVM [8], [9]), their practical performance may still fall short of expectation due to the semantic complexity of natural languages. For instance, on most ACSA tasks of the popular SemEval benchmark, the reported top accuracy levels are only around 80 percent [1], [10].

It can be observed that natural languages provide rich linguistic hints potentially useful for polarity reasoning. A sentence may contain strong sentiment words that explicitly express sentiment. In the running example, the presence of the strong sentiment word “like”, together with the absence of any negative word, suggests that the sentiment of the sentence S_{11} is positive. A sentence may also contain shift words (e.g., *but* and *however*), which do not directly indicate polarity but explicitly specify the relationship between two neighboring aspect polarities. Again in the running example, the word “However” at the beginning of the sentence S_{12} indicates that its polarity is opposite to the polarity of the sentence S_{11} . In contrast, the absence of any shift word between two neighboring sentences usually means that their polarities are similar (e.g., S_{21} and S_{22}).

Unfortunately, the existing DNN techniques have limited capability in modeling various linguistic hints. In this paper, we propose a novel framework, SenHint, which enables joint inference based on both DNN and linguistic hints. It first extracts explicit linguistic hints and then encodes their

• The authors are with the School of Computer Science, Northwestern Polytechnical University, and also with the Key Laboratory of Big Data Storage and Management, Northwestern Polytechnical University, Ministry of Industry and Information Technology, Xi’an, ShaanXi 710072, China.
E-mail: {wangyanyan, a.murtadha}@mail.nwpu.edu.cn, {chenbenben, lizhh, panwei1002, liuhailong}@nwpu.edu.cn.

Manuscript received 14 Apr. 2019; revised 19 Sept. 2019; accepted 30 Sept. 2019. Date of publication 0 . 0000; date of current version 0 . 0000.

(Corresponding author: Qun Chen.)

Recommended for acceptance by Benjamin C. M. Fung.

Digital Object Identifier no. 10.1109/TKDE.2019.2947587

TABLE 1
A Running Example From Laptop Reviews

R_i	S_{ij}	Text
R_1	S_{11} S_{12}	I like the battery that can last long time. However, the keyboard sits a little far back for me.
R_2	S_{21} S_{22}	The laptop has a long battery life. It also can run my games smoothly .

69 implications as well as the output of DNN in a unified
70 model based on Markov logic network (MLN) [11]. We note
71 that it is not new to leverage linguistic hints for sentiment
72 analysis. The traditional lexicon-based approaches [12] used
73 the hints of sentiment words to directly predict polarity by
74 summing up all the sentiment scores; the hints of context-
75 sensitive sentiment words have been integrated into deep
76 neural networks for sentiment analysis [13], [14]; the hints
77 of shift words have also been used to tune the performance
78 of deep neural networks for sentence-level sentiment analy-
79 sis [15]. However, *SenHint* is novel in that it models both the
80 output of deep neural networks and the implications of lingu-
81 guistic hints as first-class citizens in a unified MLN. Com-
82 pared with previous work, *SenHint* also leverages linguistic
83 hints for new purposes. For instance, it uses the hints of shift
84 words to capture the implicit relations between aspect
85 polarities for MLN reasoning.

86 The major contributions of this paper can be summarized
87 as follows:

- 88 1) We propose *SenHint*, a joint inference framework for
89 aspect-level sentiment analysis based on MLN. *SenHint*
90 can seamlessly integrate the output of DNN and the impli-
91 cations of linguistic hints in a unified model;
- 92 2) We present the required techniques for linguistic
93 hint extraction, MLN model construction, and joint
94 MLN inference;
- 95 3) We empirically evaluate the performance of *SenHint*
96 on both English and Chinese benchmark datasets. Our
97 extensive experiments show that compared to the state-of-
98 the-art DNN techniques, *SenHint* can effectively im-
99 prove polarity detection accuracy by considerable margins.

100 Note that a prototype of *SenHint* has been demonstrated
101 in [16]. We summarize the new contributions of this techni-
102 cal paper as follows:

- 103 1) It proposes an improved MLN model. Besides the
104 implicit polarity relations indicated by the presence/
105 absence of shift words, the new MLN model also en-
106 codes the influence of sentiment words on aspect
107 polarities;
- 108 2) It presents the improved techniques for linguistic
109 hint extraction, MLN model construction, and joint
110 inference. Unlike the demo paper, it provides with
111 the technical details of each proposed technique;
- 112 3) In empirical evaluation, while the demo paper only
113 applied *SenHint* to ACSA tasks, it extends *SenHint*
114 to handle both ACSA and ATSA tasks. Besides the
115 DNN models used in the demo paper, it also com-
116 pares *SenHint* to the more recently proposed DNN
117

118 techniques for both ACSA and ATSA. It also sepa-
119 rately evaluates the effect of various linguistic hints
120 on the performance of *SenHint*. Finally, it empiri-
121 cally compares the new *SenHint* with the original
122 version proposed in the demo paper. The experi-
123 ments have shown that the new *SenHint* performs
124 evidently better.

125 The rest of this paper is organized as follows: Section 2
126 reviews more related work. Section 3 defines the task and
127 introduces Markov logic network, the reasoning model
128 underlying *SenHint*. Section 4 gives the overview of the pro-
129 posed framework. Section 5 presents the techniques of
130 extracting linguistic hints. Section 6 describes how to encode
131 the implications of linguistic hints as well as the output of
132 DNN in a MLN. Section 7 presents the technique of joint
133 inference. Section 8 presents the empirical evaluation results.
134 Finally, we conclude this paper with some thoughts on
135 future work in Section 9. 136

2 RELATED WORK 137

138 In general, sentiment analysis involves various tasks, such as
139 polarity classification, subjectivity or objectivity identifica-
140 tion, and multimodal fusion [17]. In this paper, we focus on
141 the essential task of polarity classification. Sentiment analy-
142 sis at different granularity levels, including document, sen-
143 tence, and aspect levels, has been extensively studied in the
144 literature [18].

145 *Document and Sentence Level Sentiment Analysis.* At the doc-
146 ument (resp. sentence) level, its goal is to detect the polarity
147 of the entire document (resp. sentence) without regard to the
148 mentioned aspects. The state-of-the-art approaches were
149 built on deep neural networks (e.g., CNN and RNN), which
150 include Character-level Convolutional Networks [19], Deep
151 Pyramid Convolutional Neural Networks [20] and Linguisti-
152 cally Regularized LSTM [14]. Many works proposed to com-
153 bine an attention mechanism with neural networks, for
154 instance Hierarchical Attention Network [21], Hierarchical
155 Query-driven Attention Network [22], Linguistic-aware
156 Attention Network [23] and Cognition Based Attention
157 Model [24]. Moreover, Self-Attention Network [25] (inspired
158 by the Transformer architecture), a flexible and interpretable
159 architecture, has been proposed for text classification. Unfor-
160 tunately, all these proposals can not be directly applied to
161 aspect-level sentiment analysis because a sentence may hold
162 different opinions on different aspects.

163 *Aspect-Level Sentiment Analysis.* Aspect-level sentiment
164 analysis needs to first extract the target aspects from a given
165 sentence, and then determine their sentiment polarities. The
166 popular models for aspect extraction, which include Atten-
167 tion Based Aspect Extraction [26] and Aspect Extraction
168 with Sememe Attentions [27], employed unsupervised
169 framework analogous to an autoencoder to learn the aspects
170 with various attention mechanisms. There also exist some
171 work aiming to jointly detect the aspects and identify their
172 sentiment polarity [28], [29].

173 In this paper, we instead focus on how to determine the
174 polarities of the given aspects in a sentence. Since deep neu-
175 ral networks can automatically learn high-quality features or
176 representations, the state-of-the-art approaches attempted to
177 adapt such models for aspect-level sentiment analysis. The

TABLE 2
Frequently Used Notations

Notation	Description
$t_i = (r_j, s_k, a_l)$	an aspect unit
r_j	a review
s_k	a sentence
a_l	an aspect category or aspect term
$T = \{t_i\}$	a set of aspect units
$v(t_i)$	a boolean variable indicating whether the sentiment polarity of t_i is positive
$V = \{v(t_i)\}$	a set of aspect polarity variables

existing work can be divided into two categories based on the two finer subtasks of ATSA and ACSA.

For ATSA task, Dong [30] initially proposed an Adaptive Recursive Neural Network (AdaRNN) that can employ a novel multi-compositionality layer to propagate the sentiments of words towards the target. Noticing that the models based on recursive neural network heavily rely on external syntactic parser, which may result in inferior performance, the following work [31] focused on recurrent neural networks. The alternative solutions include memory networks [32] and convolutional neural networks [33]. Due to the great success of attention mechanism in machine translation [34] and question answering [35], many models based on LSTM and attention mechanism have also been proposed. These models, including Hierarchical Attention Network [36], Segmentation Attention Network [37], Interactive Attention Networks [38], Recurrent Attention Network [39], Attention-over-Attention Neural Networks [40], Effective Attention Modeling [41], Content Attention Model [42], Multi-grained Attention Network [43], employed different attention mechanisms to output the aspect-specific sentiment features. More recently, the capsule networks [44], a type of artificial neural network that can better model hierarchical relationships, have also been leveraged for ATSA task. Chen [45] proposed a Transfer Capsule Network for transferring document-level knowledge to aspect-level sentiment analysis.

In comparison, there exist fewer works for ACSA because the implicit aspects make the task more challenging. Ruder [46] proposed a hierarchical bidirectional LSTM that can model the inter-dependencies of sentences in a review. Wang [47] presented an attention-based LSTM that employs an aspect-to-sentence attention mechanism to concentrate on the key part of a sentence given an aspect. Xue [3] introduced a model based on convolutional neural networks and gating mechanisms. Wang [48] presented an AS-Capsule model that can fully employ the correlation between aspect and sentiment through shared components. Note that the models proposed for ACSA can also be used for ATSA, but the ones for ATSA usually solely benefit themselves because they usually employ specific components to model an explicit aspect term together with its relative context.

Other Relevant Work. Word representation, which has been used as input by all the DNN models, plays an important role in sentiment analysis. Traditional word representations [49] are effective at capturing semantic and syntactic information, but they usually perform poorly in capturing sentiment polarity. Therefore, there exist some work on sentiment-specific work representation. For instance,

Tang [50], [51] proposed C&W based models to learn sentiment-specific word embedding by distant supervision for twitter sentiment classification. Fu [52] employed local context information as well as global sentiment representation to learn sentiment-specific word embeddings.

Markov logic network, as an expressive template language, enables joint inference based on both feature and relational information. It has been widely applied to many applications [11]. However, the existing approaches based on MLN generally require human-designed features. In this paper, we integrate the DNN output and linguistic hints into a unified model based on MLN, which can retain the relational reasoning ability of MLN while avoiding complicated feature engineering.

3 PRELIMINARIES

In this section, we first define the task and then introduce Markov logic network (MLN), the inference model underlying SenHint.

3.1 Task Statement

For presentation simplicity, we have summarized the frequently used notations in Table 2. We formulate the task of aspect-level sentiment analysis as follows:

Definition 1 [Aspect-level Sentiment Analysis]. Let $t_i = (r_j, s_k, a_l)$ be an aspect unit, where r_j is a review, s_k is a sentence in the review r_j , and a_l is an aspect associated with the sentence s_k . Note that the aspect a_l can be a aspect category or aspect term, and a sentence may express opinions towards multiple aspects. Given a corpus of reviews, R , the goal of the task is to predict the sentiment polarity of each aspect unit t_i in R .

3.2 Markov Logic Network

Markov logic network combines first-order logic and probabilistic graphical model in a single representation. In first-order logic, a set of formulas represent the hard constraints over a set of instances, and the rules can not be violated. The basic idea of MLN is to generalize first-order logic by softening the hard constraints, assigning a weight to each formula to indicate its strength. In MLN, the instances can violate the formulas but need to pay a penalty: the higher the weight, the greater the penalty to be paid. Formally, a MLN is defined as follows:

Definition 2 [Markov Logic Network]. A MLN consists of a collection of weighted first-order logic formulas $\{(F_i, w_i)\}$, where F_i is a formula in first-order logic and w_i is a real number indicating the level of confidence on this formula.

An example of MLN has been shown in Table 3.

Grounding. A MLN provides a template for constructing factor graph. A factor graph consists of variable vertices $X = \{x_1, \dots, x_n\}$ and factor vertices $\Phi = \{\phi_1, \dots, \phi_n\}$, where each factor ϕ_i is a function $\phi_i(X_i)$ over the variables X_i ($X_i \subset X$). The factors together define a joint probability distribution over all the variables X .

Provided with a MLN and a set of constants, the process of constructing factor graph is called *grounding* [53]. In the grounding process, for each predicate and formula in MLN, we will create a set of *ground atoms* and *ground formulas*, which

TABLE 3
An Example of MLN and its Corresponding Predicates and Constants

Weight	First-order logic	Predicate	Person(P)	Fact
2.0	$smoke(x) \rightarrow cancer(x)$	$smoke(x)(x \in P)$	Anna	friend(Anna, Bob)
3.0	$smoke(x) \wedge friend(x, y) \rightarrow smoke(y)$	$cancer(x)(x \in P)$	Bob	
		$friend(x, y)(x, y \in P)$		

TABLE 4
Grounding of the Example MLN (V_{id} and F_{id} Represent Variable and Factor Respectively)

V_{id}	Ground atoms	F_{id}	Ground formulas	Ground factor graph
x_1	smoke(Anna)	f_1	$smoke(Anna) \rightarrow cancer(Anna)$	
x_2	cancer(Anna)	f_2	$smoke(Bob) \rightarrow cancer(Bob)$	
x_3	smoke(Bob)	f_3	$smoke(Anna) \wedge friend(Anna, Bob) \rightarrow smoke(Bob)$	
x_4	cancer(Bob)			

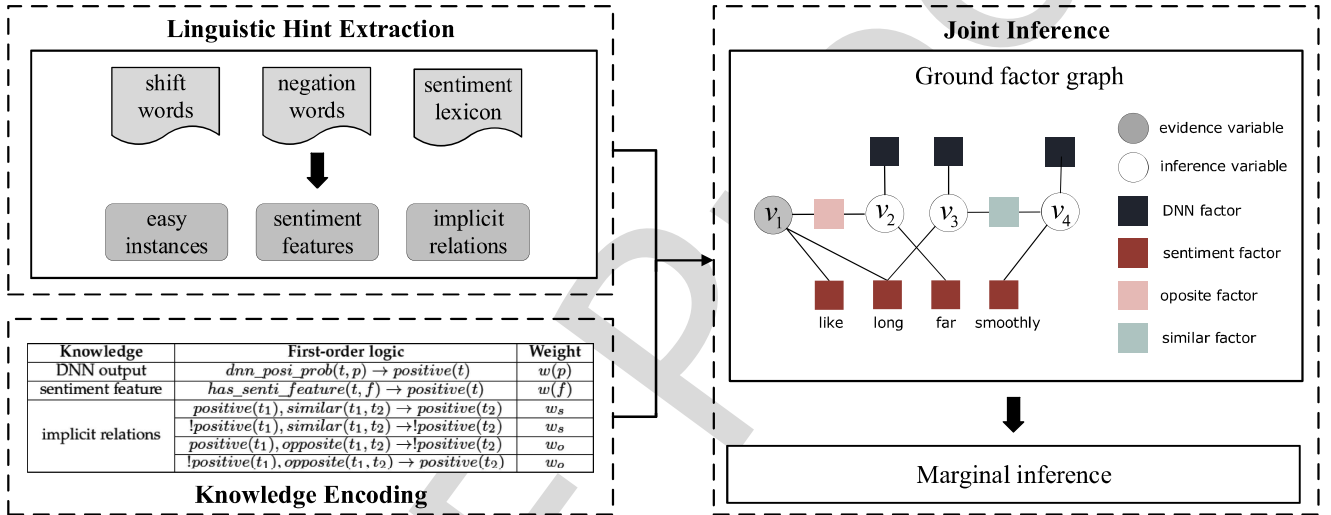


Fig. 1. The framework overview of SenHint.

are represented by the variables and factors respectively in the factor graph. The grounding process of the MLN defined in Table 3 has been shown in Table 4

Marginal Inference. A factor graph defines a joint probability distribution over its variables X by

$$P(X = x) = \frac{1}{Z} \prod_i \phi_i(X_i) = \frac{1}{Z} \exp\left(\sum_i w_i n_i(x)\right), \quad (1)$$

where n_i denotes the number of true groundings of the formula F_i in x , w_i denotes the weight of F_i , and Z is the partition function, i.e., normalization constant. The process of computing the probability of each variable is referred to as *marginal inference*.

4 FRAMEWORK OVERVIEW

As shown in Fig. 1, the framework of SenHint consists of the following three modules:

- *Linguistic Hint Extraction:* This module retrieves relevant linguistic hints from reviews. It identifies easy instances of aspect polarity, extracts common sentiment features shared by aspect polarities and mines their polarity relations.

- *Knowledge Encoding:* This module employs weighted first-order logic rules to encode the implications of linguistic hints as well as the outputs of DNN into a MLN. The outputs of DNN capture the implicit influence resulting from multiple levels of automatically learned features, while the implications of linguistic hints enable explicit polarity inference.
- *Joint Inference:* This module constructs a ground factor graph based on the generated weighted first-order logic rules, and then performs joint inference on the factor graph.

The example factor graph constructed for the running example has been shown in Fig. 1, in which aspect polarities are represented by variables (round nodes in the figure), and the influence of DNN output and linguistic implications are represented by factors (box nodes in the figure). There are two types of variables: *evidence variable* and *inference variable*. The evidence variables represent the easy instances, whose sentiment polarities can be directly determined by explicit linguistic hints with high accuracy. They participate in the inference process, but their values are specified beforehand and remain unchanged throughout the whole process. The inference variables represent the more challenging instances. Their values should instead be inferred based on the factor graph.

325 Additionally, there are four types of factors: *DNN factor*,
326 *sentiment factor*, *similar factor* and *opposite factor*. The DNN
327 factor simulates the effect of DNN output on polarity. The
328 sentiment factor captures the influence of sentiment fea-
329 tures. The similar factor and opposite factor encode the rela-
330 tions between aspect polarities.

331 5 LINGUISTIC HINT EXTRACTION

332 In this section, we describe how to identify easy instances,
333 extract sentiment features and mine polarity relations by
334 linguistic hints.

335 5.1 Identifying Easy Instances

336 The existing lexicon-based approaches essentially reason
337 about polarity by summing up the polarity scores of the sen-
338 timent words in a sentence. However, they are prone to error
339 under some ambiguous circumstances. First, the presence of
340 contrast (e.g., *but* and *although*), hypothetical (e.g., *if*) or con-
341 dition (e.g., *unless*) connectives could significantly compli-
342 cate polarity detection. For instance, the sentence “would be
343 a very nice laptop if the mousepad worked properly”
344 contains only the positive sentiment words “nice” and
345 “properly”, but it holds negative attitude due to the presence
346 of the hypothetical connective “if”. Second, the presence of
347 negation words involving long-distance dependency could
348 also make the task challenging. For instance, in the sentence
349 “I don’t really think the laptop has a good battery life”, the
350 negation word “don’t” reverses the polarity, but it is far
351 away from the sentiment word “good”. Unfortunately, the
352 existing approaches for negation detection based on local
353 neighborhood [12] can not work properly in the circum-
354 stance of long-distance dependency. Finally, a sentence may
355 not contain strong sentiment words, or even if it does, multi-
356 ple sentiment words may hold conflicting polarities. For
357 instance, consider the sentence “To be honest, i am a little
358 disappointed and considering returning it”. Since it contains
359 both the positive word “honest” and the negative word
360 “disappointed”, its true polarity is not easily detectable
361 based on sentiment word scoring.

362 Therefore, for easy instance identification, SenHint choo-
363 ses to exclude the instances with the aforementioned ambig-
364 uous patterns. Specifically,

365 **Definition 3 [Easy Instances].** *SenHint identifies an aspect*
366 *polarity as an easy instance if and only if the sentence express-*
367 *ing opinions about the aspect satisfies the following three*
368 *conditions:*

- 369 • *It contains at least one strong sentiment word, but does*
370 *not simultaneously contain any sentiment word hold-*
371 *ing the conflicting polarity;*
- 372 • *It does not contain any contrast, hypothetical or condi-*
373 *tion connective;*
- 374 • *It does not contain any negation word involving long-*
375 *distance dependency;*

376 In SenHint, the polarity of an easy instance is simply
377 determined by the polarity of its strong sentiment word. Sen-
378 Hint considers a sentiment word as *strong* if and only if the
379 absolute value of its score exceeds a pre-specified threshold
380 (e.g., 1.0 in our experiment, where the scores of sentiment

words are normalized into the interval of [-4,4]). Moreover, a
381 negation word is supposed to involve long-distance depen-
382 dency if and only if it is not in the neighboring 3-grams pre-
383 ceding any sentiment word. We illustrate the difference
384 between the easy and challenging instances by Example 1. 385

Example 1 [Easy Instances]. In a phone review, the sen- 386
387 tence “the screen is not good for carrying around in your
388 bare hands”, which expresses the opinion about “screen”,
389 is an easy instance, because the sentiment word “good”
390 associated with the local negation cue “not” strongly indi-
391 cates the negative sentiment. In contrast, the sentence “I
392 don’t know why anyone would want to write a great
393 review about this battery”, which expresses the opinion
394 about “battery”, is not an easy instance. Even though it
395 contains the strong sentiment word “great”, it includes the
396 negation word “don’t” involving long-distance depen-
397 dency. Similarly, the sentence “I like this laptop, the only
398 problem is that it can not last long time” is not an easy
399 instance, because it contains both the positive and negative
400 words (e.g., “like” and “problem”).

401 5.2 Extracting Sentiment Features

402 Sentiment words usually play an important role in deter-
403 mining the aspect polarities in a sentence. Accordingly, two
404 sentences sharing a sentiment word usually have the same
405 sentiment polarity. Hence, SenHint extracts the common
406 sentiment words from sentences and model their influence
407 by feature factors in the unified MLN model. Sentiment fea-
408 tures include both the generic sentiment words in an open-
409 source lexicon developed by Liu [2], or the domain-specific
410 sentiment words¹ that can be automatically mined from the
411 unlabeled review corpora. Since negation words can effec-
412 tively reverse polarity, we also perform negation detection
413 for each sentiment word by examining whether there is any
414 negation in its neighboring words.

415 To enable more accurate influence modeling, we also
416 propose to filter sentiment features based on the syntactic
417 structure of sentence. First, SenHint uses the constituency
418 based parse tree [54] to identify sentence structure (e.g.,
419 compound or complex) and then determines the important
420 part of a sentence based on the structure. Specifically, if a
421 sentence describes only one aspect and has a compound
422 structure with the coordinating conjunction “but”, we only
423 retain the sentiment features appearing in the “but” clause.
424 Second, in the case that multiple aspects are opined in a sen-
425 tence, SenHint uses the dependency based parse tree [55] to
426 extract the opinion phrases, each of which is a pair of opin-
427 ion target and word, for the mapping between the sentiment
428 features and their target aspects. Specifically, it associates an
429 opinion word (corresponding to a sentiment feature) with
430 an aspect if and only if either its opinion target or the opin-
431 ion word itself is close to the aspect term in the vector space.
432 We illustrate sentiment feature extraction by Example 2.

Example 2 [Sentiment Feature Extraction]. Consider the 433
434 sentence, “I thought learning the Mac OS would be hard,
435 but it is easily picked up”, which expresses the opinion
436 about the aspect “os#usability”. SenHint extracts “easily”

1. <http://www.wowbigdata.cn/SenHint/SenHint.html>

437 as sentiment feature but not “hard”, because the word
 438 “hard” does not appear in the “but” clause. Consider
 439 another example, “The screen is gorgeous, and the perfor-
 440 mance is excellent.”, which comments on both aspects
 441 of “display#quality” and “laptop#performance”. SenHint
 442 extracts two opinion phrases $\langle screen, gorgeous \rangle$ and
 443 $\langle performance, excellent \rangle$, and then reasons that 1)
 444 “gorgeous” is a feature of the aspect “display#quality”
 445 because its opinion target “screen” is very close to the
 446 aspect in vector space; 2) “excellent” is a feature of the
 447 aspect “laptop#performance” because the aspect term
 448 explicitly appears in the opinion phrase.

449 5.3 Mining Polarity Relations

450 Modeling sentences independently, the existing DNNs for
 451 aspect-level sentiment analysis have very limited capability
 452 in capturing contextual information at sentence level. How-
 453 ever, sentences build upon each other. There often exist some
 454 discourse relations between sentences that can provide valu-
 455 able hints for sentiment prediction [56]. The most influential
 456 discourse relation is the contrast relation, which is often
 457 marked by shift words (e.g., *but* and *however*). Specifically,
 458 two sentences connected with a shift word usually have oppo-
 459 site polarities. In contrast, two neighboring sentences without
 460 any shift word between them usually have similar polarities.

461 Based on these observations, SenHint extracts the similar
 462 and opposite relations between aspect polarities based on
 463 sentence context. Given two aspect units $t_i = \{r_i, s_i, a_i\}$ and
 464 $t_j = \{r_j, s_j, a_j\}$ that occur in the same review (namely
 465 $r_i = r_j$), the rules for extracting polarity relations are defined
 466 as follows:

- 467 1) If the sentences s_i and s_j are identical ($s_i = s_j$) or adja-
 468 cent and neither of them contains any shift word, t_i
 469 and t_j are supposed to hold similar polarities;
- 470 2) If two adjacent sentences s_i and s_j are connected by a
 471 shift word and neither of them contains any inner-
 472 sentence shift word, t_i and t_j are supposed to hold
 473 opposite polarities;
- 474 3) If the sentences s_i and s_j are identical and the opin-
 475 ion clauses associated with them are connected by a
 476 inner-sentence shift word, t_i and t_j are supposed to
 477 hold opposite polarities.

478 We illustrate polarity relation mining by Example 3.

479 **Example 3 [Polarity Relation Mining].** In the running
 480 example shown in Table 1, the aspect polarities in S_{21} and
 481 S_{22} are supposed to be similar based on the 1st rule. Since
 482 S_{11} and S_{12} in R_1 are connected by the shift word of
 483 “However”, their aspect polarities are reasoned to be
 484 opposite based on the 2nd rule. Additionally, consider the
 485 sentence “The screen is bright but the processing power is
 486 not very good”, which expresses the opinions about both
 487 “screen” and “processing power”. It can be observed that
 488 the two opinion clauses are connected by the shift word
 489 “but” within the sentence. Therefore, their polarities are
 490 supposed to be opposite based on the 3rd rule.

491 6 KNOWLEDGE ENCODING IN MLN

492 Note that SenHint models the easy instances of aspect polar-
 493 ity as evidence variables in MLN. In this section, we describe

494 how to encode the output of DNN, sentiment features and
 495 polarity relations in MLN.

496 6.1 Encoding DNN Output

497 In this paper, we use the recently proposed gated convo-
 498 lutional networks [3] (GCAE) as an illustrative example.
 499 The outputs of other DNNs can however be encoded in
 500 SenHint in the same way. GCAE uses convolutional neu-
 501 ral networks and gating mechanisms to selectively output
 502 the sentiment features associated with a given aspect. Its
 503 output can indicate the influence resulting from multiple
 504 levels of features that correspond to different levels of
 505 abstraction.

506 SenHint encodes the influence of DNN outputs using the
 507 following rule:

$$508 w(p) : dnn_posi_prob(t, p) \rightarrow positive(t), \quad (2)$$

509 in which $dnn_posi_prob(t, p)$ predicates that the probability
 510 of an aspect unit t having the positive polarity is equal to
 511 the value of p , $positive(t)$ is a boolean variable indicating the
 512 polarity of t , and $w(p)$ denotes the level of confidence on the
 513 rule. Observing that the relationship between the weight w
 514 and the probability p (for a boolean variable x being true)
 515 can be expressed by $p(x = 1) = e^w / (1 + e^w)$, we define the
 516 rule weight as
 517

$$518 w(p) = \ln\left(\frac{p}{1-p}\right). \quad (3)$$

519 According to Eq. (3), $w(p) > 0$ if $p > 0.5$; otherwise, if
 520 $p < 0.5$, then $w(p) < 0$. In the case of $w(p) > 0$, a zero
 521 value of $positive(t)$ would invoke a cost penalty as desired.
 522 In the case of $w(p) < 0$, a positive value for $positive(t)$
 523 would instead invoke a cost penalty.
 524
 525

526 6.2 Encoding Sentiment Features

527 SenHint encodes the influence of sentiment features using
 528 the following rule:

$$529 w(f) : has_senti_feature(t, f) \rightarrow positive(t), \quad (4)$$

530 where $has_senti_feature(t, f)$ predicates that the aspect unit
 531 t has the sentiment feature f , and $w(f)$ denotes the feature
 532 weight. In our implementation, the weight of a sentiment
 533 feature is initially set to 1 if it is a positive word in the lexi-
 534 con, or -1 if it is a negative word. Based on the labeled
 535 instances, SenHint learns the weights of sentiment features
 536 in joint inference, and their learned values are supposed to
 537 reflect their sentiment intensity. For instance, in the factor
 538 graph as shown in Fig. 1, the variable v_1 contains two senti-
 539 ment features “like” and “long”, and the sentiment feature
 540 of “long” is also shared by v_3 . Both sentiment features
 541 have positive weights, and the learned weight of “like”
 542 holds a higher value than the learned weight of “long”.
 543 Their weights accurately reflect their relative sentiment
 544 intensity.
 545

546 6.3 Encoding Polarity Relations

547 SenHint encodes the influence of similar relation between
 548 two aspect polarities by

$$w_s : \text{positive}(t_1), \text{similar}(t_1, t_2) \rightarrow \text{positive}(t_2), \quad (5)$$

and

$$w_s : \neg \text{positive}(t_1), \text{similar}(t_1, t_2) \rightarrow \neg \text{positive}(t_2), \quad (6)$$

in which w_s denotes a positive weight, t_1 and t_2 denote two aspect units and $\neg \text{positive}(t_i)$ denotes the negation of a boolean variable. For instance, in the factor graph as shown in Fig. 1, there exists a similar relation between v_3 and v_4 , which represent the instances in S_{21} and S_{22} respectively. As expected, the encoding rules of Eqs. (5) and (6) would force them to hold similar polarity, otherwise a cost penalty would be invoked.

Similarly, SenHint encodes the influence of opposite relation between two aspect polarities by

$$w_o : \text{positive}(t_1), \text{opposite}(t_1, t_2) \rightarrow \neg \text{positive}(t_2), \quad (7)$$

and

$$w_o : \neg \text{positive}(t_1), \text{opposite}(t_1, t_2) \rightarrow \text{positive}(t_2), \quad (8)$$

in which w_o denotes a positive weight.

SenHint interprets rule weight or confidence on rule as the accuracy of mined relations. With the polarity of t_1 being positive, the probability of the polarity of t_2 being positive can be estimated by

$$p(v(t_2) = 1) = e^{w_s} / (1 + e^{w_s}). \quad (9)$$

Approximating $p(v(t_2) = 1)$ with the accuracy r_{acc} , we can establish the relationship between rule weight and relation accuracy by

$$w_s = \ln \left(\frac{r_{acc}}{1 - r_{acc}} \right). \quad (10)$$

SenHint sets the rule weight w_o specified in (7) and (8) in a similar way. Note that the higher the estimated accuracy, the higher the rule weights. For accuracy estimation, SenHint first applies the mining rules to the labeled data used for DNN training, and then approximates the accuracy on the test data with the result observed on the training data. Our empirical evaluation in Section 8.3 has shown that the accuracies achieved on the test data are generally high, and very similar to the results observed on the training data in most cases.

7 JOINT INFERENCE

The MLN model of SenHint is comprised of the formulas specified in Eqs. (2), (4), (5), (6), (7) and (8). Based on the model, SenHint first constructs a factor graph, and then estimates the marginal probabilities of inference variables.

Denoting the DNN, sentiment, similar, opposite factors by $\phi_p^{dnn}(\cdot)$, $\phi_f^{sent}(\cdot)$, $\phi^{sim}(\cdot, \cdot)$, $\phi^{opp}(\cdot, \cdot)$ respectively, SenHint defines them as follows:

$$\phi_p^{dnn}(v(t)) = \begin{cases} 1 & v(t) = 0, \\ e^{w(p)} & v(t) = 1. \end{cases} \quad (11)$$

$$\phi_f^{sent}(v(t)) = \begin{cases} 1 & v(t) = 0, \\ e^{w(f)} & v(t) = 1. \end{cases} \quad (12)$$

$$\phi^{sim}(v(t_1), v(t_2)) = \begin{cases} 1 & v(t_1) \neq v(t_2), \\ e^{w_s} & v(t_1) = v(t_2). \end{cases} \quad (13)$$

$$\phi^{opp}(v(t_1), v(t_2)) = \begin{cases} 1 & v(t_1) \neq v(t_2), \\ e^{-w_o} & v(t_1) = v(t_2). \end{cases} \quad (14)$$

where $v(t)$ denotes a boolean variable indicating the polarity of t , and $w(p)$, $w(f)$, w_s and w_o denote the rule weights.

Based on the factors, the factor graph defines a joint probability distribution over its variables V by

$$P_w(V) = \frac{1}{Z} \prod_{v \in V} \phi_p^{dnn}(v(t)) \prod_{v \in V} \prod_{f \in F_v} \phi_f^{sent}(v(t)) \prod_{(t_1, t_2) \in R} \phi^{rel_type}(v(t_1), v(t_2)), \quad (15)$$

where F_v denotes the set of sentiment features associated with the variable v , R denotes the sets of polarity relations between aspect units, rel_type denotes the relation type of the aspect units t_1 and t_2 (namely *sim* or *opp*) and Z denotes a partition function, i.e., normalization constant.

Given a factor graph with some labeled evidence variables, SenHint reasons about the factor weights by minimizing the negative log marginal likelihood as follows:

$$\hat{w} = \arg \min_w -\log \sum_{V_I} P_w(\Lambda, V_I), \quad (16)$$

where Λ denotes the observed labels of evidence variables and V_I denotes the set of inference variables. The objective function effectively learns the factor weights most consistent with the label observations of the evidence variables. SenHint optimizes the objective function by leveraging the Snorkel engine [57], which interleaves stochastic gradient descent steps with Gibbs sampling ones. It has been shown in [57], [58] that similar to contrastive divergence [59], the optimization process can guarantee convergence. For more details, please refer to the literature of [57], [58]. Note that in our implementation, the weights $w(p)$, w_s , w_o are automatically set to be fixed values based on the formulas of Eqs. (3) and (10), while the weight $w(f)$ is learned by optimizing the objective function. Once the weights are learned, SenHint performs marginal inference over the factor graph to compute the probability distribution for each inference variable $v(t) \in V$. SenHint uses the Numbskull library² for marginal inference.

8 EMPIRICAL EVALUATION

In this section, we empirically evaluate the performance of SenHint on the benchmark datasets by a comparative study. We compare SenHint with the state-of-the-art DNN models proposed for ACSA and ATSA. For the ACSA tasks, the compared models include:

- *H-LSTM* [46]. The hierarchical bidirectional LSTM can model the inter-dependencies of sentences in a review;
- *AT-LSTM* [47]. The Attention-based LSTM (AT-LSTM) employs an attention mechanism to concentrate on the key parts of a sentence given an aspect, where the

2. <https://github.com/HazyResearch/numbskull>

TABLE 5
Details of Benchmark Datasets

Data	Language	Train		Test	
		#T(ACSA)	#T(ATSA)	#T(ACSA)	#T(ATSA)
PHO16	Chinese	1333	—	529	—
CAM16	Chinese	1259	—	481	—
LAP16	English	2715	1478	751	435
RES16	English	2134	1662	693	578
LAP15	English	1864	1049	868	410
RES15	English	1410	1154	725	508

aspect embeddings are used to determine the attention weight;

- *ATAE-LSTM* [47]. The Attention-based LSTM with Aspect Embedding (ATAE-LSTM) extends AT-LSTM by appending the input aspect embedding into each word input vector;
- *GCAE* [3]. The gated convolutional network employs CNN and gating mechanisms to selectively output the sentiment features according to a given aspect.

For the ATSA, the compared models include:

- *IAN* [38]. The interactive attention network interactively learns the attentions in the contexts and targets, and generates the representations for targets and contexts separately;
- *RAM* [39]. The multiple-attention network can effectively capture sentiment features separated by a long distance, and is usually more robust against irrelevant information;
- *AOA* [40]. The attention-over-attention network models aspects and sentences in a joint way, and can explicitly capture the interaction between aspects and context sentences;
- *TNet* [33]. Compared with previous alternatives, the target-specific transformation network can better integrate target information into the word representations.

The rest of this section is organized as follows: Section 8.1 describes the experimental setup. Section 8.2 presents the comparative evaluation results. Section 8.3 separately evaluates the effect of easy instances, sentiment features and aspect polarity relations on the performance of SenHint. Finally, Section 8.4 presents the results of error analysis on SenHint for its future improvement.

8.1 Experimental Setup

We used the benchmark datasets in four domains (phone, camera, laptop and restaurant) and two languages (Chinese and English) from the SemEval 2015 task 12 [10] and 2016 task 5 [1]. Our experiments performed 2-class classification to label an aspect polarity as *positive* or *negative*, and thus ignored the neutral instances in our experiments. The statistics of the test datasets are presented in Table 5, in which *PHO*, *CAM*, *LAP* and *RES* denote the domain phone, camera, laptop and restaurant respectively, and #T(ACSA) and #T(ATSA) denote the numbers of aspect category units and aspect term units respectively. Since there are no labeled aspect terms in the Chinese datasets, we compare SenHint to its alternatives only on the English datasets for ATSA. Note that given a test dataset, the number of instances in its

factor graph is equal to the number of aspect category units or aspect term units it contains.

In our experiments, we used the GCAE model to predict the DNN output, because it has been empirically shown to outperform other DNN alternatives. However, SenHint can easily integrate any other DNN model into its MLN. For identifying easy instances, we used the Opinion Lexicon³ and EmotionOntology⁴ lexicons for English and Chinese data respectively. Due to their limited numbers, we manually specified the negation and shift words. In the implementation of SenHint joint inference, the number of learning and inference epochs is set at 1,000, the step size for learning is set at 0.01, the decay for updating step size is set at 0.95, and the regularization penalty is set at $1e - 6$. More details on the experimental setup can be found in our technical report [60]. Our implementation codes have also been made open-source.⁵

8.2 Comparative Evaluation

We have compared performance on both metrics of accuracy and macro-F1. Note that the metric of macro-F1 is the unweighted average of the F1-score for each label. The comparative results on the ACSA and ATSA tasks are presented in Tables 6 and 7 respectively, in which *SenHint(demo)* denotes the original model presented in our demo paper [16] and *SenHint* denotes the improved model proposed in this paper. We have highlighted the best performance on each test task by *bold* in the tables. It can be observed that for ACSA, SenHint achieves better performance than the DNN approaches on all the test datasets. It achieves the improvement of more than 4 percent on 5 out of totally 6 tasks (i.e., PHO16, CAM16, LAP16, LAP15 and RES15). For ATSA, the experimental results are similar. SenHint outperforms the best DNN model by around 7 percent on LAP15 and LAP16, and by around 4 percent on RES15. Due to the widely recognized challenge of sentiment analysis, the achieved improvements can be considered to be very considerable. These experimental results clearly demonstrate the efficacy of SenHint.

It is also worthy to point out that *SenHint* consistently performs better than *SenHint(demo)*. The achieved improvements on most tasks are between 1 and 3 percent. The maximal improvement of around 3.5 percent is achieved on the LAP16 workload of ATSA. The only exception is PHO16, on which *SenHint* performs slightly worse than *SenHint(demo)* by less than 0.1 percent if measured by macro-F1. Our experimental results have evidently validated the efficacy of the improved MLN model proposed in this paper.

To further validate the efficacy of extracted linguistic hints, we have also conducted ablation test on both ACSA and ATSA tasks. The evaluation results have been shown in Tables 6 and 7, where *SenHint(w/o easy)*, *SenHint(w/o sentiments)* and *SenHint(w/o relations)* denote the ablated models with the components of easy instances, sentiment features and polarity relations being removed from SenHint respectively. It can be observed that: 1) SenHint achieves better performance than the ablated models in most cases with only a few exceptions. It means that all the extracted

3. <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html~liub/FBS/sentiment-analysis.html>

4. <http://ir.dlut.edu.cn/EmotionOntologyDownload>

5. <http://www.wowbigdata.cn/SenHint/SenHint.html>

TABLE 6
Performance Comparison for ACSA on Benchmark Datasets

Model	PHO16		CAM16		LAP16		RES16		LAP15		RES15	
	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1
H-LSTM	73.30%	72.59%	78.80%	73.04%	78.90%	77.18%	83.10%	79.48%	80.00%	78.25%	77.10%	76.15%
AT-LSTM	72.40%	72.16%	81.70%	77.42%	76.03%	74.73%	85.03%	80.57%	81.03%	79.10%	77.25%	77.00%
ATAE-LSTM	74.48%	73.85%	83.36%	79.59%	79.07%	77.10%	84.66%	80.50%	80.68%	78.97%	79.13%	77.83%
GCAE	76.03%	75.49%	82.49%	76.72%	80.75%	79.24%	86.87%	83.07%	81.96%	80.56%	81.49%	80.45%
SenHint(demo)	80.45%	80.20%	86.58%	82.89%	83.07%	81.71%	88.09%	84.73%	84.60%	83.46%	82.50%	81.78%
SenHint(w/o easy)	80.72%	80.08%	87.82%	84.29%	85.57%	84.26%	89.32%	86.01%	87.28%	86.20%	85.24%	84.58%
SenHint(w/o senti-feats)	80.08%	79.53%	87.53%	83.87%	84.69%	83.28%	89.00%	85.73%	86.84%	85.75%	85.43%	84.84%
SenHint(w/o relations)	80.00%	79.40%	87.82%	84.37%	82.61%	81.24%	87.07%	83.40%	86.08%	85.01%	83.83%	83.06%
SenHint	80.89%	80.15%	88.10%	84.47%	85.60%	84.28%	89.09%	85.72%	87.46%	86.40%	85.84%	85.34%

TABLE 7
Performance Comparison for ATSA on Benchmark Datasets

Model	LAP16		RES16		LAP15		RES15	
	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1	Acc	Macro-F1
AT-LSTM	74.85%	72.39%	84.43%	77.50%	77.51%	74.41%	75.43%	71.57%
ATAE-LSTM	75.08%	71.93%	84.60%	76.82%	77.66%	73.83%	74.13%	69.67%
GCAE	78.34%	75.74%	88.86%	81.93%	81.37%	79.08%	77.60%	71.81%
IAN	74.02%	71.90%	85.12%	77.01%	79.27%	76.30%	75.00%	69.34%
RAM	77.47%	75.33%	85.81%	78.44%	78.58%	76.33%	73.23%	66.33%
AOA	74.94%	72.27%	87.02%	75.83%	80.73%	77.84%	73.43%	69.71%
TNet	75.86%	73.85%	87.20%	80.20%	80.00%	78.88%	75.20%	71.32%
SenHint(demo)	82.75%	80.98%	89.65%	83.25%	86.47%	84.75%	81.17%	77.53%
SenHint(w/o easy)	85.47%	83.82%	89.79%	84.08%	87.90%	86.28%	80.87%	76.73%
SenHint(w/o senti-feats)	84.78%	83.22%	89.69%	84.03%	87.66%	86.10%	81.77%	78.10%
SenHint(w/o relations)	84.32%	82.53%	88.93%	82.91%	87.27%	85.66%	81.02%	77.02%
SenHint	86.19%	84.65%	89.68%	84.12%	87.98%	86.41%	81.66%	77.98%

linguistic hints are helpful for polarity reasoning; 2) Among the ablated models, SenHint(w/o relations) achieves the overall worst performance, followed by SenHint(w/o senti-feats) and SenHint(w/o easy). It means that the influence of polarity relations on the performance of SenHint is the greatest, followed by sentiment features and easy instances.

It can also be observed that the improvement margins of SenHint over *SenHint(w/o easy)* and *SenHint(w/o senti-feats)* are very similar on the English and Chinese datasets; however, the influence of polarity relations is greater on the English datasets than the Chinese datasets. In the experiments, we have observed that more polarity relations can be extracted from the English datasets than the Chinese datasets, and they are generally accurate. Therefore, as shown in Table 6, *SenHint* outperforms the ablated model of *SenHint(w/o relations)* by more considerable margins on the English datasets than the Chinese datasets.

8.3 Separate Effect Evaluation

In this subsection, we report our evaluation results on the ACSA tasks. The evaluation results on the ATSA tasks are similar, thus omitted here due to space limit. But they can be found in our technical report [60].

Easy Instances. We first evaluate the performance of the technique proposed for identifying easy instances. We compare its performance with the best DNN model of GCAE. Note that SenHint identifies easy instances by pre-specified rules. Therefore, for SenHint, the percentage of easy instances, which is calculated by dividing the number of easy instances by the total number of instances in a test dataset, is fixed for each test dataset. For fair comparison, we also select the same number of least uncertain instances in a test dataset based on the output of GCAE, and then compare the achieved accuracy of SenHint and GCAE. The detailed results on the ACSA tasks are presented in Table 8, in which

TABLE 8
Performance Evaluation of Identifying Easy Instances

	ACSA					
	PHO16	CAM16	LAP16	RES16	LAP15	RES15
Prop	35.73%	43.87%	46.34%	55.70%	54.72%	47.17%
Acc(GCAE)	86.35%	87.49%	90.80%	92.75%	88.76%	88.54%
Acc(SenHint)	95.24%	98.58%	93.68%	93.01%	95.16%	93.57%

TABLE 9
Performance Comparison Between GCAE and SenHint-Easy

	ACSA					
	PHO16	CAM16	LAP16	RES16	LAP15	RES15
GCAE	76.03%	82.49%	80.75%	86.87%	81.96%	81.49%
SenHint-easy	79.23%	87.32%	82.13%	86.97%	85.50%	83.82%

TABLE 10
Performance Evaluation of Polarity Relation Mining

Relation type	Data type	ACSA					
		PHO16	CAM16	LAP16	RES16	LAP15	RES15
similar relations	train	89.39%	88.89%	92.57%	95.12%	93.39%	96.07%
	test	85.71%	92.13%	93.38%	95.34%	90.51%	92.53%
opposite relations	train	75.00%	89.29%	83.33%	72.22%	80.00%	75.00%
	test	100%	90.00%	50.00%	66.67%	100%	60.00%

792 the first row (*Prop*) denotes the percentage of easy instances
793 identified by SenHint, and the following two rows (*Acc*)
794 denote the accuracy of GCAE and SenHint respectively. It
795 can be observed that

- 796 1) A considerable percentage of the instances in a test
797 workload can be identified as easy instances by Sen-
798 Hint: the percentage varies from 35 to 55 percent;
- 799 2) SenHint detects the polarities of easy instances with
800 the consistently higher accuracy than GCAE, and the
801 improvement margins are considerable. On PHO16
802 and CAM16, the margins are as large as 9-10 percent;

803 We then evaluate the effect of identified easy instances on
804 the performance of SenHint by comparing SenHint-easy
805 with GCAE, in which SenHint-easy represents the MLN
806 model using the outputs of DNN and easy instances but
807 not mined sentiment features and polarity relations. The
808 detailed results are presented in Table 9. It can be observed
809 that the MLN model of using easy instances alone can effec-
810 tively improve the performance of polarity classification. On
811 the difference between the English and Chinese datasets, we
812 have observed that a higher percentage of instances can be
813 identified as easy on the English datasets, but the achieved
814 accuracy is generally lower. Overall, their effect on the per-
815 formance of SenHint are quite similar on the English and
816 Chinese datasets.

817 *Polarity Relations.* We first evaluate the performance of
818 the technique proposed for mining polarity relations. The
819 detailed results are presented in Table 10, which reports the
820 accuracy of mined relations on both training and test data.
821 As expected, the achieved accuracies on the test data are
822 generally similar to the results obtained on the training
823 data. Most importantly, the accuracy of mined relations is
824 high ($\geq 80\%$) in most cases.

825 We then compare SenHint-rel with GCAE, in which
826 SenHint-rel denotes the MLN model integrating DNN out-
827 puts and mined polarity relations but not easy instances and
828 sentiment features. The comparative results are presented in
829 Table 11. It can be observed that SenHint-rel can effectively
830 improve the performance of DNN. These observations vali-
831 date the effectiveness of the proposed strategy, which
832 assigns different weights to relations such that a relation

with higher accuracy can have greater impact on its con- 833
834 nected variables.

835 *Sentiment Features.* We evaluate the effect of extracted senti- 835
836 ment features on the performance of SenHint by comparing 836
837 GCAE with SenHint-sent, in which SenHint-sent denotes the 837
838 MLN model integrating DNN output and extracted senti- 838
839 ment features but not easy instances and mined polarity rela- 839
840 tions. Their comparative results are presented in Table 12. 840
841 We can observe that SenHint-sent can effectively improve 841
842 the performance of DNN. These experiments validate the 842
843 effectiveness of the proposed strategy for integrating com- 843
844 mon sentiment features into the MLN model. 844

845 8.4 Error Analysis

846 For the improvement of SenHint in the future, it is helpful to 846
847 scrutinize its failure cases. We have categorized the failure 847
848 cases into the following categories: 848

- 849 • *Lack of linguistic hints.* This type of error occurs when 849
850 no linguistic hint has been extracted from a sentence. 850
851 If an instance does not any extracted linguistic hint, 851
852 its predicted polarity is the same as the DNN output. 852
853 For instance, consider the single sentence in a review, 853
854 "I would have kept it but that was the sole reason for 854
855 my purchase" , which expresses the opinion about 855
856 "laptop#general". It contains neither sentiment fea- 856
857 ture nor polarity relation. Since it is mislabeled by 857
858 DNN, SenHint also fails. 858
- 859 • *Incorrect linguistic hints.* This type of error occurs when 859
860 the extracted linguistic hints are incorrect. Most of the 860
861 errors under this category can be further categorized 861
862 into the following two subcategories: 1) the instances 862
863 are incorrectly identified as easy; 2) the extracted 863
864 polarity relations are erroneous. For the first subcate- 864
865 gory, consider the sentence, "I have to clean it regu- 865
866 larly for it to stay looking good". SenHint identifies it 866
867 as an easy instance with the positive polarity. How- 867
868 ever, its true polarity is negative. For the second sub- 868
869 category, consider two neighboring sentences, "it 869
870 looks sleek ad gorgeous" and "i find myself adjusting 870
871 it regularly". Since they are not connected by any shift 871
872 word, SenHint reasons that their polarities are similar. 872

TABLE 11
Performance Comparison Between GCAE and SenHint-Rel

	ACSA					
	PHO16	CAM16	LAP16	RES16	LAP15	RES15
GCAE	76.03%	82.49%	80.75%	86.87%	81.96%	81.49%
SenHint-rel	76.88%	82.58%	83.70%	90.93%	84.72%	82.33%

TABLE 12
Performance Comparison Between GCAE and SenHint-Sent

	ACSA					
	PHO16	CAM16	LAP16	RES16	LAP15	RES15
GCAE	76.03%	82.49%	80.75%	86.87%	81.96%	81.49%
SenHint-sent	78.26%	85.25%	81.67%	87.39%	84.09%	82.00%

TABLE 13
Distribution of Classification Errors

No.	Error category	Percentage
1	Lack of linguistic hints	32.11%
2	Incorrect linguistic hints	30.28%
3	Ineffectual linguistic hints	25.69%
4	Others	11.92%

However, they are indeed opposite. SenHint first identifies the polarity of the first sentence as positive and then incorrectly labels the polarity of the second sentence as positive based on the extracted polarity relation.

- *Ineffectual linguistic hints.* In this case, even though the extracted linguistic hints are correct, they fail to correct the erroneous outputs of DNN. For instance, consider two neighboring instances with the same positive polarity. Even though SenHint correctly extracts the similar polarity relation between them, it may still fails under the following two circumstances: 1) DNN erroneously labels both instances as negative. Since the erroneous outputs of DNN happen to satisfy the supposed relation, SenHint can not flip their polarities; 2) DNN correctly identifies one of them as positive with a lower confidence (e.g., 0.6) while erroneously identifying the other one as negative with a higher confidence (e.g., 0.05). Instead of correcting the error of DNN, SenHint may flip the polarity of the correctly identified instance from positive to negative.

Using the ACSA task on LAP16 as the test case, we have given the relative percentages of different error classes in Table 13. It can be observed that the error class of Lack of Linguistic Hints occupies the largest portion, followed by Incorrect Linguistic Hints, which comes second. Thus, improving the accuracy and coverage of linguistic hints extraction may greatly enhance the performance of SenHint.

9 CONCLUSION

In this paper, we have proposed the SenHint framework for aspect-level sentiment analysis that can integrate deep neural networks and linguistic hints in a coherent MLN inference model. We have presented the required techniques for extracting linguistic hints, encoding their implications into the model, and joint inference. Our extensive experiments on the benchmark data have also validated its efficacy.

Built on DNN, SenHint still requires considerable training data. It is interesting to observe that provided with sufficient review corpus, employing easy instance detection, extracted sentiment features and polarity relations can potentially make it unnecessary to classify aspect polarity by DNN. In future work, we will explore how to make SenHint perform well while requiring little or even no labeled training data.

ACKNOWLEDGMENTS

This work is supported by the Ministry of Science and Technology of China, National Key Research and Development Program (2018YFB1003403), National Natural Science

Foundation of China (61732014 and 61672432), and Natural Science Basic Research Plan in Shaanxi Province of China (2018JM6086).

REFERENCES

- [1] M. Pontiki et al., "Semeval-2016 task 5: Aspect based sentiment analysis," in *Proc. 10th Int. Workshop Semantic Eval.*, 2016, pp. 19–30.
- [2] B. Liu, *Sentiment Analysis and Opinion Mining*. San Rafael, CA, USA: Morgan & Claypool, 2012.
- [3] W. Xue and T. Li, "Aspect based sentiment analysis with gated convolutional networks," in *Proc. 56th Annu. Meet. Assoc. Comput. Linguistics*, 2018, pp. 2514–2523.
- [4] H. H. Do, P. W. C. Prasad, A. Maag, and A. Alsadoon, "Deep learning for aspect-based sentiment analysis: A comparative review," *Expert Syst. Appl.*, vol. 118, pp. 272–299, 2019.
- [5] K. Schouten and F. Frasincar, "Survey on aspect-level sentiment analysis," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 3, pp. 813–830, Mar. 2016.
- [6] X. Ding, B. Liu, and P. S. Yu, "A holistic lexicon-based approach to opinion mining," in *Proc. Int. Conf. Web Search Web Data Mining*, 2008, pp. 231–240.
- [7] M. Taboada, J. Brooke, M. Tofiloski, K. D. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Comput. Linguistics*, vol. 37, no. 2, pp. 267–307, 2011.
- [8] T. Mullen and N. Collier, "Sentiment analysis using support vector machines with diverse information sources," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2004, pp. 412–418.
- [9] S. Kiritchenko, X. Zhu, C. Cherry, and S. Mohammad, "NRC-Canada-2014: Detecting aspects and sentiment in customer reviews," in *Proc. 8th Int. Workshop Semantic Eval.*, 2014, pp. 437–442.
- [10] M. Pontiki, D. Galanis, H. Papageorgiou, S. Manandhar, and I. Androutsopoulos, "SemEval-2015 task 12: Aspect based sentiment analysis," in *Proc. 9th Int. Workshop Semantic Eval.*, 2015, pp. 486–495.
- [11] P. M. Domingos and D. Lowd, *Markov Logic: An Interface Layer for Artificial Intelligence*. San Rafael, CA, USA: Morgan & Claypool, 2009.
- [12] C. J. Hutto and E. Gilbert, "VADER: A parsimonious rule-based model for sentiment analysis of social media text," in *Proc. 8th Int. Conf. Weblogs Social Media*, 2014, pp. 216–225.
- [13] Z. Teng, D. Vo, and Y. Zhang, "Context-sensitive lexicon features for neural sentiment analysis," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2016, pp. 1629–1638.
- [14] Q. Qian, M. Huang, J. Lei, and X. Zhu, "Linguistically regularized LSTM for sentiment classification," in *Proc. 55th Annu. Meet. Assoc. Comput. Linguistics*, 2017, pp. 1679–1689.
- [15] Z. Hu, X. Ma, Z. Liu, E. H. Hovy, and E. P. Xing, "Harnessing deep neural networks with logic rules," in *Proc. 54th Annu. Meet. Assoc. Comput. Linguistics*, 2016, pp. 2410–2420.
- [16] Y. Wang et al., "SenHint: A joint framework for aspect-level sentiment analysis by deep neural networks and linguistic hints," in *Demonstrations Proc. Web Conf.*, 2018, pp. 207–210.
- [17] E. Cambria, "Affective computing and sentiment analysis," *IEEE Intell. Syst.*, vol. 31, no. 2, pp. 102–107, Mar./Apr. 2016.
- [18] K. Ravi and V. Ravi, "A survey on opinion mining and sentiment analysis: Tasks, approaches and applications," *Knowl.-Based Syst.*, vol. 89, pp. 14–46, 2015.
- [19] X. Zhang, J. J. Zhao, and Y. LeCun, "Character-level convolutional networks for text classification," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2015, pp. 649–657.
- [20] R. Johnson and T. Zhang, "Deep pyramid convolutional neural networks for text categorization," in *Proc. 55th Annu. Meet. Assoc. Comput. Linguistics*, 2017, pp. 562–570.
- [21] Z. Yang, D. Yang, C. Dyer, X. He, A. J. Smola, and E. H. Hovy, "Hierarchical attention networks for document classification," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics: Human Lang. Technol.*, 2016, pp. 1480–1489.
- [22] L. Luo et al., "Beyond polarity: Interpretable financial sentiment analysis with hierarchical query-driven attention," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 4244–4250.
- [23] Z. Lei, Y. Yang, and Y. Liu, "LAAN: A linguistic-aware attention network for sentiment analysis," in *Proc. Web Conf.*, 2018, pp. 47–48.
- [24] Y. Long, L. Qin, R. Xiang, M. Li, and C. Huang, "A cognition based attention model for sentiment analysis," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2017, pp. 462–471.

- 995 [25] G. Letarte, F. Paradis, P. Giguère, and F. Laviolette, "Importance
996 of self-attention for sentiment analysis," in *Proc. Workshop: Analyzing*
997 *Interpreting Neural Netw. NLP*, 2018, pp. 267–275.
- 998 [26] R. He, W. S. Lee, H. T. Ng, and D. Dahlmeier, "An unsupervised
999 neural attention model for aspect extraction," in *Proc. 55th Annu.*
1000 *Meet. Assoc. Comput. Linguistics*, 2017, pp. 388–397.
- 1001 [27] L. Luo et al., "Unsupervised neural aspect extraction with Sememes,"
1002 in *Proc. 28th Int. Joint Conf. Artif. Intell.*, 2019, pp. 5123–5129.
- 1003 [28] M. Saeidi, G. Bouchard, M. Liakata, and S. Riedel, "SentiHood:
1004 Targeted aspect based sentiment analysis dataset for urban
1005 neighbourhoods," in *Proc. 26th Int. Conf. Comput. Linguistics: Tech.*
1006 *Papers*, 2016, pp. 1546–1556.
- 1007 [29] Y. Ma, H. Peng, and E. Cambria, "Targeted aspect-based sentiment
1008 analysis via embedding commonsense knowledge into an attentive
1009 LSTM," in *Proc. 32nd Conf. Artif. Intell.*, 2018, pp. 5876–5883.
- 1010 [30] L. Dong, F. Wei, C. Tan, D. Tang, M. Zhou, and K. Xu, "Adaptive
1011 recursive neural network for target-dependent Twitter sentiment
1012 classification," in *Proc. 52nd Annu. Meet. Assoc. Comput. Linguistics*,
1013 2014, pp. 49–54.
- 1014 [31] D. Tang, B. Qin, X. Feng, and T. Liu, "Effective LSTMs for target-
1015 dependent sentiment classification," in *Proc. 26th Int. Conf. Comput.*
1016 *Linguistics*, 2016, pp. 3298–3307.
- 1017 [32] S. Wang, S. Mazumder, B. Liu, M. Zhou, and Y. Chang, "Target-
1018 sensitive memory networks for aspect sentiment classification," in
1019 *Proc. 56th Annu. Meet. Assoc. Comput. Linguistics*, 2018, pp. 957–967.
- 1020 [33] X. Li, L. Bing, W. Lam, and B. Shi, "Transformation networks for
1021 target-oriented sentiment classification," in *Proc. 56th Annu. Meet.*
1022 *Assoc. Comput. Linguistics*, 2018, pp. 946–956.
- 1023 [34] O. Firat, K. Cho, and Y. Bengio, "Multi-way, multilingual neural
1024 machine translation with a shared attention mechanism," in *Proc.*
1025 *Conf. North Amer. Chapter Assoc. Comput. Linguistics: Human Lang.*
1026 *Technol.*, 2016, pp. 866–875.
- 1027 [35] X. He and D. Golub, "Character-level question answering with
1028 attention," in *Proc. Conf. Empirical Methods Natural Lang. Process.*,
1029 2016, pp. 1598–1607.
- 1030 [36] J. Cheng, S. Zhao, J. Zhang, I. King, X. Zhang, and H. Wang, "Aspect-
1031 level sentiment classification with HEAT (hierarchical attention)
1032 network," in *Proc. Conf. Inf. Knowl. Manage.*, 2017, pp. 97–106.
- 1033 [37] B. Wang and W. Lu, "Learning latent opinions for aspect-level
1034 sentiment classification," in *Proc. 32nd AAAI Conf. Artif. Intell.*, 2018,
1035 pp. 5537–5544.
- 1036 [38] D. Ma, S. Li, X. Zhang, and H. Wang, "Interactive attention net-
1037 works for aspect-level sentiment classification," in *Proc. 26th Int.*
1038 *Joint Conf. Artif. Intell.*, 2017, pp. 4068–4074.
- 1039 [39] P. Chen, Z. Sun, L. Bing, and W. Yang, "Recurrent attention network
1040 on memory for aspect sentiment analysis," in *Proc. Conf. Empirical*
1041 *Methods Natural Lang. Process.*, 2017, pp. 452–461.
- 1042 [40] B. Huang, Y. Ou, and K. M. Carley, "Aspect level sentiment classi-
1043 fication with attention-over-attention neural networks," in *Proc.*
1044 *11th Int. Conf. Social Cultural Behavioral Model.*, 2018, pp. 197–206.
- 1045 [41] R. He, W. S. Lee, H. T. Ng, and D. Dahlmeier, "Effective attention
1046 modeling for aspect-level sentiment classification," in *Proc. 27th*
1047 *Int. Conf. Comput. Linguistics*, 2018, pp. 1121–1131.
- 1048 [42] Q. Liu, H. Zhang, Y. Zeng, Z. Huang, and Z. Wu, "Content atten-
1049 tion model for aspect based sentiment analysis," in *Proc. Conf.*
1050 *World Wide Web*, 2018, pp. 1023–1032.
- 1051 [43] F. Fan, Y. Feng, and D. Zhao, "Multi-grained attention network for
1052 aspect-level sentiment classification," in *Proc. Conf. Empirical*
1053 *Methods Natural Lang. Process.*, 2018, pp. 3433–3442.
- 1054 [44] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between
1055 capsules," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2017,
1056 pp. 3856–3866.
- 1057 [45] Z. Chen and T. Qian, "Transfer capsule network for aspect level
1058 sentiment classification," in *Proc. 57th Conf. Assoc. Comput. Linguis-*
1059 *tics*, 2019, pp. 547–556.
- 1060 [46] S. Ruder, P. Ghaffari, and J. G. Breslin, "A hierarchical model of
1061 reviews for aspect-based sentiment analysis," in *Proc. Conf. Empir-*
1062 *ical Methods Natural Lang. Process.*, 2016, pp. 999–1005.
- 1063 [47] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM
1064 for aspect-level sentiment classification," in *Proc. Conf. Empirical*
1065 *Methods Natural Lang. Process.*, 2016, pp. 606–615.
- 1066 [48] Y. Wang, A. Sun, M. Huang, and X. Zhu, "Aspect-level senti-
1067 ment analysis using AS-capsules," in *Proc. Web Conf.*, 2019,
1068 pp. 2033–2044.
- 1069 [49] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors
1070 for word representation," in *Proc. Conf. Empirical Methods Natural*
1071 *Lang. Process.*, 2014, pp. 1532–1543.
- 1072 [50] D. Tang, F. Wei, N. Yang, M. Zhou, T. Liu, and B. Qin, "Learning
1073 sentiment-specific word embedding for twitter sentiment classi-
1074 fication," in *Proc. 52nd Annu. Meet. Assoc. Comput. Linguistics*,
1075 2014, pp. 1555–1565.
- 1076 [51] D. Tang, F. Wei, B. Qin, N. Yang, T. Liu, and M. Zhou, "Sentiment
1077 embeddings with applications to sentiment analysis," *IEEE Trans.*
1078 *Knowl. Data Eng.*, vol. 28, no. 2, pp. 496–509, Feb. 2016.
- 1079 [52] P. Fu, Z. Lin, F. Yuan, W. Wang, and D. Meng, "Learning sentiment-
1080 specific word embedding via global sentiment representation," in
1081 *Proc. 32nd AAAI Conf. Artif. Intell.*, 2018, pp. 4808–4815.
- 1082 [53] Y. Chen and D. Z. Wang, "Knowledge expansion over probabilistic
1083 knowledge bases," in *Proc. ACM SIGMOD Int. Conf. Manage.*
1084 *Data*, 2014, pp. 649–660.
- 1085 [54] M. P. Marcus, B. Santorini, and M. A. Marcinkiewicz, "Building a
1086 large annotated corpus of English: The penn treebank," *Comput.*
1087 *Linguistics*, vol. 19, no. 2, pp. 313–330, 1993.
- 1088 [55] D. Chen and C. D. Manning, "A fast and accurate dependency
1089 parser using neural networks," in *Proc. Conf. Empirical Methods*
1090 *Natural Lang. Process.*, 2014, pp. 740–750.
- 1091 [56] E. Pitler, M. Raghupathy, H. Mehta, A. Nenkova, A. Lee, and
1092 A. K. Joshi, "Easily identifiable discourse relations," in *Proc.*
1093 *22nd Int. Conf. Comput. Linguistics*, 2008, pp. 87–90.
- 1094 [57] A. Ratner, S. H. Bach, H. R. Ehrenberg, J. A. Fries, S. Wu, and
1095 C. Ré, "Snorkel: Rapid training data creation with weak super-
1096 vision," *Proc. VLDB Endowment*, vol. 11, no. 3, pp. 269–282, 2017.
- 1097 [58] S. H. Bach, B. D. He, A. Ratner, and C. Ré, "Learning the structure
1098 of generative models without labeled data," in *Proc. 34th Int. Conf.*
1099 *Mach. Learn.*, 2017, pp. 273–282.
- 1100 [59] G. E. Hinton, "Training products of experts by minimizing contrastive
1101 divergence," *Neural Comput.*, vol. 14, no. 8, pp. 1771–1800, 2002.
- 1102 [60] Y. Wang, Q. Chen, M. Ahmed, Z. Li, W. Pan, and H. Liu, "Joint infer-
1103 ence for aspect-level sentiment analysis by deep neural networks
1104 and linguistic hints(technical report)," 2019. [Online]. Available:
1105 <http://www.wowbigdata.com.cn/SenHint/senhint-report.pdf>



Yanyan Wang is working toward the PhD degree
with the School of Computer Science, Northwestern
Polytechnical University. Her research interests
include sentiment analysis and artificial intelligence.



Qun Chen is a professor with the School of
Computer Science, Northwestern Polytechnical
University. His current research interests include
interdisciplinary methodologies and techniques
(mostly based on data analysis and machine learn-
ing) for a variety of challenging computation tasks
(e.g., entity resolution and sentiment analysis).



Murtadha Ahmed is working toward the PhD
degree with the School of Computer Science, North-
western Polytechnical University. His research
interests include sentiment analysis and artificial
intelligence.

1122
1123
1124
1125
1126
1127



Zhanhuai Li is a professor with the School of Computer Science, Northwestern Polytechnical University. His research interests include data storage and management. He has served as program committee chair or member in various conferences and committees.



Hailong Liu is an associate professor with the School of Computer Science, Northwestern Polytechnical University. His research interests include data quality management.

1132
1133
1134
1135

1128
1129
1130
1131



Wei Pan is an associate professor with the School of Computer Science, Northwestern Polytechnical University. His research interests include graph processing.

▷ **For more information on this or any other computing topic,** please visit our Digital Library at www.computer.org/publications/dlib.

1136
1137

IEEE Proof