

Evaluative Pattern Extraction for Automated Text Generation

Chia-Chen Lee and Shu-Kai Hsieh

Graduate Institute of Linguistics, National Taiwan University
she767219@gmail.com, shukai@gmail.com

Abstract

Getting travel tips from the experienced bloggers and online forums has been one of the important supplements to the travel guidebook in the web society. In this paper we present a novel approach by identifying and extracting evaluative patterns, providing a different linguistically-motivated framework for automated evaluative text generation. We target at domain-specific observation in online travel blogs in Chinese. Results suggest that the semantic prosody accompanying the patterns demonstrates that online travel bloggers prefer to employ tacit pragmatic strategy in presenting their sentiment polarity in comments. The extracted patterns and their differentiation can be beneficial to identifying and characterizing evaluative language for further automated opinion summarization and macro/micro planning in natural language generation (NLG) as well.

1 Introduction

With the rapidly growing use of the Internet, text mining, sentiment analysis, and evaluative language analysis of online resources are becoming essential issues. Online travel blogs serve as main opinions and comments providers sharing their traveling experiences where the texts are constructed with authors' evaluation about the traveling. The automation of text planning in this domain has become highly demanded. This paper aims to propose a linguistic framework of working with evaluative expressions by examining domain-restricted specialized discourse of traveling articles. Identifying the particular linguistic behaviors and patterns of evaluative language agglomerative structure would facilitate both the macro/ micro planning in NLG in this domain.

In online travel blog articles, evaluative language is expressed in several kinds. lexical level terms

such as 'recommend', 'delicious', and 'surprise', are explicit evaluations. Other than this, patterns are found and can be generalized into a certain fixed meanings in traveling domain. For instance, 有N味 'has the flavor/feeling of N' is a common pattern used as in 有家鄉味 'has the feeling of home', 有台灣味 'has the feeling of Taiwan' as positive evaluation in the data. We propose to adopt pattern grammar (Hunston, 1999) in approaching the evaluative prosody widely occurred in the travel blogs. Pattern grammar focuses on the concept that meaning belongs to patterns, targeting on the recurring co-occurrences and the particularly shared meanings of lexical item nodes. There is specialized domain-specific grammar not applying to general grammar, resulting a fixed meaning of patterns in that specific domain. As Sinclair (1991) said: "It seems that there is a strong tendency for sense and syntax to be associated", suggesting that meaning and its patterns are highly related. Francis (1993) used the pattern *v it adj* as an example, which limits the choices of its lexical items on either verbs or adjectives, indicating that the meaning of a pattern is also limited and patterns will occur with words through semantic restriction. Therefore, patterns extracted from the texts should be the primary consideration and observation for natural language processing, particularly for semantic and sentiment analysis, whether as for annotation, summarization or text generation.

2 Literature Review

In NLG, content determination is an essential process to decide what is the communicated information in texts (Reiter, 1995). In order to generate natural-language text, a system must be able to determine what to include and how to organize the information to achieve its communicative goal most effectively. McKeown (1985) based on discourse strategies as a guide for natural-language text generation, which generated paragraph-length

responses. In domain-specific texts such as weather forecast (Adeyanju, 2012), automated text generation is expected to have similar weather conditions where its language pattern is observable. In traveling blog articles, the evaluative language is its dominant feature. Evaluative language has been researched since 1970s, starting from Halliday (1976), with others making further developments or moving on to new approaches such as Chafe (1986), Biber and Finegan (1989), Hunston (1994), Francis (1995), and Martin and White (2000). Hunston (1994, 2000, 2004) defined evaluative language as which is “expressed through language which indexes the act of evaluation or the act of stance-taking. It expresses an attitude towards a person, situation, or other entity and is both subjective and located within a societal value system”. It is the driving force behind virtually all communications. (Thompson and Hunston, 2000). Patterns of a word are defined as “all the words and structures which are regularly associated with the word and which contribute to its meaning”. The relationship between patterns and lexis is mutually dependent, in that each pattern appears with a limited set of lexical items, and each lexical item occurs with a restricted set of patterns. As patterns are highly associated with meaning, words sharing a given pattern will also tend to share an aspect of meaning (Hunston, 1999).

With the concepts combination of evaluative language and pattern grammar, we can discover that how evaluation is spread across texts with fixed meanings. The necessity of examining evaluation language is obvious in that online travel blog articles serve as the purpose for sharing comments and opinions to readers, and to find out if there are certain structures or patterns in the texts are utilizable for generating opinion summaries.

3 Patterns and Evaluative Meanings in Content Determination

The categorization of evaluation languages is diverse for different research purposes. To fit the communicative goal in the traveling context, where *recommendation* instead of neutral descriptions is needed, the following relevant attributes are targeted: attraction, hotel, restaurant, food, and event. Among these targets, evaluative expressions are realized in different aspects. For instance, main

evaluated aspects for attraction are its environment, transportation, popularity, culture, and so on. While in food, its price, taste, quality, or quantity are main discussed issues. **Table 1** shows the attributes and their evaluated aspects.

Attributes	Evaluated Aspects
Attraction	Environment (space, design, atmosphere, weather), transportation, popularity
Hotel	Environment (space, design, atmosphere), transportation, popularity, price, service
Restaurant	Environment (space, design, atmosphere), transportation, popularity, price, service
Food	Popularity, price, taste, quality, quantity
Event	Environment (space, design, atmosphere, weather), popularity, product(price, package, quality)

Table 1: Lightweight ontology in traveling domain

In this study, data are crawled from ten online travel blogs nominated as the ten most popular online travel blogs in GOLDDOT Award 2015¹, held by Pixnet in Taiwan, with 540 articles in total. A corpus-based approach is taken for exploring the data and extracting the patterns. As evaluated patterns are embodied within sentences and flexible in its unit, there is no straightforward way to observe them in the corpus. Annotation is based on the attributes mentioned earlier for categorization, using LOPOTATOR, an online linguistic annotation tool designed by LOPE lab². One annotator is involved in annotation process. Chunks are considered as units for patterns detection, mostly restricted in phrasal units, where the evaluator and the evaluation are included so as to know the relationship between the property of evaluated entity and the evaluation expression. For instance, chunk like 值得一探的美景 ‘a beautiful view that is worth visiting’ will be annotated as with the evaluator 美景 ‘beautiful view’ and its expression 值得一探的 ‘something which is worth visiting’. The processing pipeline is shown in **Figure 1**.

¹ <http://2015golddot.events.pixnet.net/>

² <http://lopen.linguistics.ntu.edu.tw:8001/lope.anno>

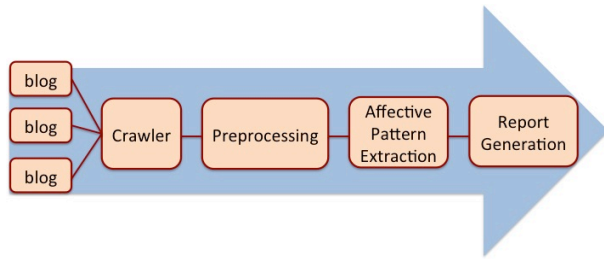


Figure 1: Processing pipeline.

4 Data Annotation and Analysis

Different from previous linguistic formalisms (such as Rhetorical Structure Theory) used in document structuring, where the main focus is hierarchical construct of messages, the *evaluative pattern grammar* as proposed in this paper explores the linear interaction of lexis and configuration at the evaluative level. In our corpus, lexical items are explicitly observable evaluation, such as 大 ‘big’, 新 ‘new’, 好 ‘good’, 分享 ‘share’, 推薦 ‘recommend’, 喜歡 ‘like’, and 享受 ‘enjoy’ are frequently occurred in the data. Our primary attention here is to extract the fixed patterns denoting fossilized polarity in evaluation co-occurring with a variety of word choices.

Manual annotation for patterns extraction in online travel blog articles provides an exhaustive result of all possible evaluative use.

In all annotated units, expressions with similar meanings and structures can be generalized into patterns, generating a fixed basic meaning, where they seem to be neutral but denote a polarity when used in a context. **Table 2** summaries the patterns listed by different aspects, with a symbol ‘+’ and ‘-’ representing the polarity being positive or negative the pattern implies. Due to limit of pages a few patterns are listed as instances. Whenever a pattern occurs, it brings out a value merging with the meaning of its variant noun, verb, or adjectives. 非常有 N 味 ‘so full of N’s flavor or feeling’ is taken as an example. In this pattern, it’s the comment on the food evaluator that it is ‘full of the flavor or feeling’ of the noun phrase, with implicit neutral evaluation until noun phrase is filled in, such as 非常有家鄉味 and realized as the meaning of ‘full of home’s feeling; the food makes you feel or think of home’, gaining positive evaluation.

Patterns	Instances	Polarity
[N 直逼 N] ‘N can nearly compete with N’	設計感 直逼 W Hotel ‘its design can compete with W Hotel’	+
[N 有梗] ‘N is interesting’	空間 有梗 ‘the space is interesting’	+
[N 破表] ‘it’s quite over of the degree of N’	浪漫指數 破表 ‘it’s quite over of the degree of the romance’	+
[讓你有種 N 的感覺] ‘make you have the feeling of N’	讓你有種 家 的感覺 ‘make you have the feeling to be home’	+
[絕對是 N 的 N] ‘it’s definite N’s N’	絕對是 飯店 的 基本配備 ‘it’s definite the basic equipment of a hotel’	+
[N 對我來說已是 N] ‘N is already N to me’	甜度 對我來說已是 極限 ‘the sweetness is already way too enough to me’	-
[非常有 N 味] ‘so full of N’s feeling’	非常有 家鄉 味 ‘so full of home’s feeling’	+
[N 十足] ‘a lot of N; high degree of N’	咬勁 十足 ‘high degree of texture’	+
[光是 V 就知它的 N] ‘knowing its N just by V’	光是 看顏色 就知它的 粉嫩程度 ‘knowing its freshness just by looking at the color’	+

Table 2: Evaluative patterns and data instances.

Patterns shown in **Table 2** are case-specific to the traveling domain, and they can be taken as self-embedded evaluative meaning carriers which are

useful cues in content determination in that a pattern can simply be a comment unit shown a posi-

tive or negative evaluation toward the evaluated targets.



Figure 2: User interface snapshots of traveling recommendation searching and searching results.

Figure 2 is a temporary template of user interface where users can search for traveling comments or opinions, and the comments can be either using the evaluative patterns generated from our work or the origin sentences from the author.

Comments from several authors' comments and scores of the traveling targets are useful when only searching for a single and specific target, such as Taipei 101 or W Hotel. However, common occasions are that people want to know all possible comments on one target, such as recommendation for traveling in Tokyo, with all things might be experienced in Tokyo. Therefore, we create a simplified plan (exemplified in English version) as in **Figure 3** for generating the evaluative summary from a single author's traveling article. Parenthesis units such as '(name of the author)' in **Figure 3** are information to be extracted from the article, including author's name, places or things experienced by the author with comments. Evaluators are comment units extracting from our pattern generation work. Both opinions are informative generation results.

The blogger *_(author's name)_* came to *_(traveling places)_* for traveling, where he/she experienced *_(place1)_*, *_(place2)_*, *_(place3)_*, and *_(place4)_*. About *_(evaluator 1)_*, *_(name of the author)_* like because he/she thinks that it is *_(evaluative pattern 1)_*, particularly *_(part of evaluator 1)_* is worth trying. In addition, he/she also went to *_(evaluator 2)_*, and he/she recommended it because of *_(evaluative pattern 2)_*. Among

that, *_(part of evaluator 2)_* is the most recommended one. ...

Figure 3. Simplified document plan.

In short, the identification of evaluative patterns in texts, as inspired by usage-based linguistic pattern grammar theory, can be utilized as a key feature for domain-specialized research on opinion mining and generation in evaluative texts.

5 Conclusion and Future Work

Due to the socio-pragmatic reasons, the evaluative patterns found in online travel blogs have their own characteristics and therefore call for more attention. On one hand, the recurrent linguistic means of evaluation as performed in texts of this genre are mostly beyond the word level; on the other hand, bloggers often tacitly organize their discourse of feelings or assessments in a relatively polite manner. It constitutes a challenge for content selection and text planning, more linguistic framework should be involved in properly tailoring the data for potential users.

The approach proposed in this paper can handle with affective contents as seen crucial in the opinionated text mining and generation, has encountered its limitation mainly related to the annotation process. Manual annotation can achieve higher accuracy in extracting possible patterns, however subjective annotation with only one annotator causes time-consuming and inefficiency problems. There are few studies relating to the evaluative language in online traveling blog domain, this paper serves as a point of departure in discovering the evaluative patterns, and as a reference for probing into other domain-specific evaluative language. Patterns extraction can be applied to other domains and the annotated data can be used for automatic pattern extraction algorithms and for text summarization in the process of document planning in NLG. For text generation, pattern is a significant feature as a representation of the sentiment or polarity toward the evaluation. Automated patterns extraction will be a valuable progress in generating evaluative text summary.

References

- Adeyanju, I. 2012. *Generating weather forecast texts with case based reasoning*. International Journal of Computer Applications, 45.
- Biber, D., and E. Finegan. 1989. *Styles of stance in English: Lexical and grammatical marking of evidentiality and affect*. Text 9.93-124. Special issue on *The pragmatics of affect*, ed. by Elinor Ochs).
- Chafe, Wallace and Nichols, Johanna. 1986. *Evidentiality: the Linguistic Coding of Epistemology*. Norwood, New Jersey: Ablex
- Francis, G. 1993. *A corpus-driven approach to grammar — principles, methods and examples*. In Baker et al. (eds), 137–156.
- Francis, G. 1995. *Corpus-driven grammar and its relevance to the learning of English in a cross-cultural situation*. In *English in Education: Multicultural perspectives*, A. Pakir (ed). Singapore: Unipress.
- Hunston, S. and Francis, G. 1999 *Pattern Grammar: a corpus-driven approach to the lexical grammar of English*. Amsterdam: Benjamins.
- Hunston, S. and Sinclair, J. 2000 ‘A local grammar of evaluation’ in Hunston and Thompson (eds.) *Evaluation in Text: authorial stance and the construction of discourse*. Oxford: Oxford University Press.
- Martin, J. R. & White, P. R. R. 2005. *The language of evaluation: appraisal in English*. Basingstoke : Palgrave Macmillan.
- Sinclair, J.M. 1991. *Corpus, Concordance, Collocation*. Oxford: OUP.
- McKeown, K. R. 1985. *Discourse strategies for generating natural-language text*. Artificial Intelligence, 27(1), 1-41.
- Reiter E. and Dale R. 1997. *Building applied natural language generation systems*. *Natural Language Engineering*, 3, pp 57-87.