# Explaining Item Ratings in Cosmetic Product Reviews

Yuuki Matsunami, Mayumi Ueda, Shinsuke Nakajima,
Takeru Hashikami, Sunao Iwasaki, John O'Donovan, and Byungkyu Kang

*Abstract*—In the cosmetics domain, many online sellers support user-provided product reviews. It has been shown that reviews have a profound effect on product conversion rates. Reviews of cosmetic products carry particular importance in purchasing decisions because of their personal nature, and particularly because of the potential for irritation with unsuitable products. In this paper, we propose a method for automatic scoring of various aspects of cosmetic item review texts based on a curated dictionary of expressions from a corpus of real world online reviews. Results and discussion of a user experiment to evaluate the approach are presented. In particular, we find that a co-occurrence approach improved coverage of reviews, and that our automated approach predicted attributes in manually annotated ground truth with an accuracy of 81%.

*Index Terms*—explanation of reviews, automatic review rating, evaluation expression dictionary, analysis of review

Fig. 1. Example of automatic scoring of various aspects of cosmetic item reviews

## I. INTRODUCTION

IN recent years, many online sellers of cosmetic products have added support for user-provided reviews. These are very helpful for consumers to decide whether to buy a commercial product, and they have been shown to have a significant impact on conversion rates. In particular, consumers make a careful choices about cosmetics since unsuitable products frequently cause skin irritations. "@cosme"[1] is very popular among Japanese young women as one of cosmetics item review sites. The site is very helpful, however, it is not easy to find truly suitable cosmetics because of the lack of explanation and granularity in user provided item ratings. As an example, there is no guarantee that a cosmetics item, which one contributor mentioned as good for dry skin, is always suitable for people who have dry skin. Since the compatibilities between skin and cosmetics items differ from one user to another, we believe it is important to identify and cluster users who share common preferences for cosmetic products and to share reviews among those niche communities. To study the proposed approach, we design and evaluate a collaborative recommender system for cosmetic products, which incorporates opinions of similar-minded users and automatically scores fine grained aspects of product reviews.

In order to develop such a review recommender system, we have to analyze review text to understand feedbacks on reviewers' experiences of cosmetic items. Actually, there is a score (as # of stars) of each review text on the conventional cosmetic review sites. However, it is mostly overall score, so that it is difficult to recognize feedbacks on reviewers' experiences of cosmetic items from it. For examples, there are "moisturizing effect", "whitening effect", "exfoliation care effect", "Hypoallergenic effect", and "Aging care effect", etc. as aspects of reviews for "face lotion". Thus, we need a scoring method of such various aspects of cosmetic item review texts to understand feedbacks on reviewers' experiences of them.

Hence, the purpose of our study is to propose a method for automatic scoring of various aspects of cosmetic item review texts based on evaluation expression dictionary. The method can realize an automatic scoring of various aspects of cosmetic item reviews which have even if no scores (see Fig.1). In this paper, we construct an evaluation expression dictionary for "face lotion", which has different feedbacks with each person, as a first step. Moreover, we discuss the adequateness of our proposed method based on an evaluation experiment for the automatic scoring method.

The rest part of the paper is organized as follows. The related work is given in section 2. Then section 3 describes the method for automatic scoring of various aspects of cosmetic item review texts based on evaluation expression dictionary. Discussions about the adequateness of our proposed method based on an evaluation experiment for the automatic scoring method are given in section 4. We conclude the paper in section 5.

## II. Related Work

There are many websites to provide reviews by the consumers. For example, Amazon.com[2] and Priceprice.com[3] are popular shopping sites over the Internet, and these sites provide reviews of their merchandise by the consumers. And "Tabelog" is also popular website in Japan. This website does not sell products, it provides restaurant's information and reviews. In addition to the algorithmic aspects, researchers have recently focused on the presentation aspects of review data[9]. Furthermore, in recent years, "@cosme" is very popular among Japanese young women. This website is a portal site for beauty and cosmetic items, and it provides various information, such as reviews and shopping information of cosmetic items. According to the report by the istyle Inc. that is a operation company of this system, at November 2015, the number of monthly page view is 280 million, the number of member is 3.5 million, and the total number of review is 1200 million[4]. From this report, many women exchange information about beauty and cosmetics through the service of @cosme. At the service of @cosme, they provide many information about cosmetic items of various cosmetic brands. Hence, users can compare cosmetic items through the various cosmetic brands. Reviews are composed of review text, scores, tag about effects, etc. Furthermore, the system has profile data that includes information about age and skin type, made by the users when they enroll as a member. Therefore, users who want to browse the reviews can search the reviews according to their own purposes, for example, reviews sorted by the scores or focused on one effect.

Along with the popularization of these review services, several researches about analysis of reviews have been conducted in the past. For example, O'Donovan et al. evaluated their AuctionRules algorithm –a dictionary-based scoring mechanism for eBay reviews of Egyptian antiques. They showed that the approach was scalable and particualrly that a small amount of domain knowledge can greatly improve prediction accuracy compared against traditional instance-based learning approaches. In our previous study, we analyze reviews of the cosmetic items[5]. In order to determine if the review is positive review or negative review, we make dictionaries for the Japanese language morphological analysis, which composed of positive expression and negative expression of cosmetic items. This previous research is aimed to develop the system to provide the reviews that take account of the user's profile, then, that system tries to retrieve information from blogs and SNS, and merge the information to the same format. Our final goal of current study is to develop a method for automatic scoring of review texts, according to various aspects of cosmetic items.

Nihongi et al. propose a method for extracting the evaluation expression from the review texts, and they develop the product retrieval system using evaluation expressions[6]. Our research focuses on the analysis of the review for cosmetic items, and we are aimed to find similar users about preferences and feelings in order to recommend truly useful reviews.

Titov et al. propose a statistical model for sentiment summarization[7]. This model is a joint model of text and aspect ratings. In order to discover the corresponding topics, this model uses aspect ratings. Therefore, this model is able to extract textual evidence from reviews without the need of annotated data.

As stated above, there are several studies to analyze reviews. However, there has been no study that tried to develop a method for automatic scoring of review texts, according to various aspects of cosmetic items.

## III. Automatic Scoring of Various Aspects of Cosmetic Item Review Texts Based on Evaluation Expression Dictionary

In this section, we describe a method for automatic scoring of various aspects of cosmetic item review texts based on evaluation expression dictionary. At firest, we describe the brief overview of our proposed method in section III-A. Section III-B explains how to construct the evaluation expression dictionary. The method for automatic scoring of various aspects of cosmetic item review texts is given in section III-C.

### A. Overview of Proposed Method

The purpose of this paper is to propose a method for automatic scoring of various aspects of cosmetic item review texts based on evaluation expression dictionary. Furthermore, our final goal is to develop a cosmetic item review recommender system which can recommend truly useful reviews for a target user. It performs by using a set of similar users who have common both preferences and feedbacks on their experiences of the cosmetic items.

In order to make a significance of our study clear, Fig.2 shows a conceptual diagram of the cosmetic item review recommender system which is our final goal. In Fig.2, numbers in blue written as (1) - (4) are corresponding to the procedure of cosmetic review automatic scoring process, and Roman alphabets in red written as (a) - (e) are corresponding to the procedure of review recommendation process. More detailed procedure of each process are shown below:

Cosmetic Review Automatic Scoring Process

(1) Construct the evaluation expression dictionary which includes pairs of evaluation expression and its score by analyzing reviews sampled from non-scored DB.
(2) Pick up reviews from non-scored DB to score them.
(3) Automatically score reviews picked up in step (2) based on the evaluation expression dictionary constructed in step (1).
(4) Put reviews scored in step (3) into scored review DB.

Review Recommendation Process

(a) User give the name of a cosmetic item that she is interested in.
(b) System Refers to "similar user extraction module" in order to extract similar users to the target user of step (a).
(c) "Similar user extraction module" obtains the information of reviews and reviewers, and identify similar users to the target user.
(d) Provide reviews of the similar users identified in step (c) to "Review recommendation module".
(e) System recommends suitable reviews to the target user.
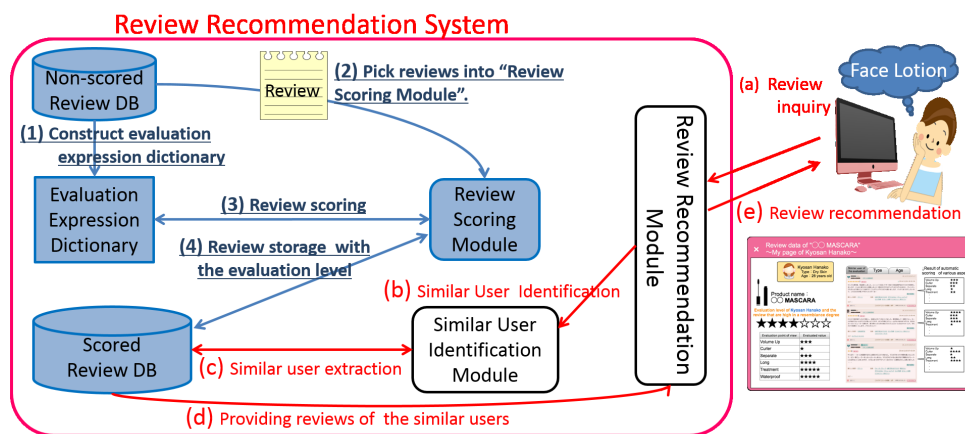
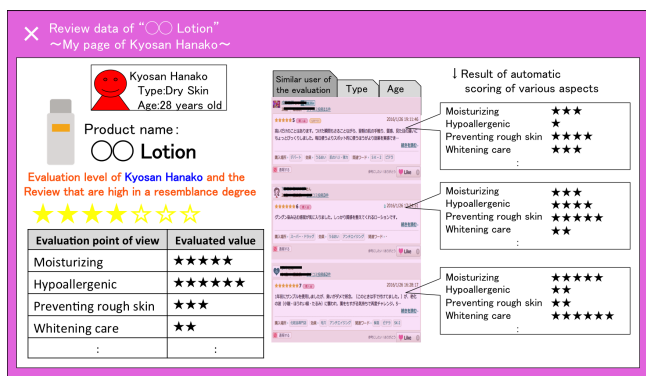Fig. 2.   Conceptual Diagram of the Cosmetic Item Review Recommender System



Fig. 3.   An example of User Interface of the Cosmetic Item Review Recommender System



Fig. 4.   Review Scoring Using Phrase Expression-based Dictionary

We try to develop the cosmetic review automatic scoring method in this paper. And developing the review recommendation method is our future work.

Fig.3 shows an example of user interface of the cosmetic item review recommender system. We believe that users can browse truly suitable reviews of a target cosmetic item and also they can easily choose reviewers group such as "reviewers having similar evaluation tastes", "reviewers having a similar skin type" and "reviewers of the same age group" by clicking the tab. Moreover, the user interface can provide not only reviews themselves but also their scores for the various aspects against the target cosmetic item, so that users can understand what kind of feedbacks on reviewers' experiences of cosmetic items without difficulty.

### B. Constructing the Evaluation Expression Dictionary

We describe how to construct the evaluation expression dictionary which has pairs of evaluation expressions and the scores against cosmetic items in this section.

Reviewers about cosmetic items post reviews against cosmetic items with widely varying expressions. Thus, we try to construct the dictionary by extracting and registering evaluation expressions from real review data. We gather review data to construct the dictionary from @cosme[1] which is the representative cosmetic review website in Japan. In particular, we extract both frequently-appearing expressions in good evaluations and bad evaluations for each cosmetic item, and register these evaluation expressions into the dictionary.

There are many kinds of cosmetic items. As a first step, we try to construct evaluation expression dictionary of "face lotion" in this paper. The reason why we focus on "face lotion" is that "face lotion" is used by a lot of people. In addition, there are various evaluations against even one product of "face lotion" due to differences of users' skin types.

*1) Phrase Expression-based Dictionary:* In order to construct the phrase expression-based dictionary, we gather 80 reviews for "face lotion", and manually extract 1,893 characteristic evaluation phrases from them. Next, two evaluators, who are 20's female students, give a score against each evaluation expression phrase manually, and we set an average of the scores as the final score of the expression phrase. There are widely various expressions in review texts because they are based on free description in natural languages. Therefore, we categorize gathered evaluation expressions into 39 groups which correspond to detailed Categories in Table I by consulting the effect-tags in @cosme. Fig.4 shows the data format of the phrase expression-based dictionary.

The procedure of automatic scoring method based on the phrase expression-based dictionary is as follows:

At first, the method gathers non-scored reviews, and identify evaluation expressions existing in these reviews. Secondly, it gives a score to each evaluation expression based on the dictionary if there is same evaluation expression in the dictionary.

For example in Fig.4, the review text includes phrases as "considerably moistened" and "moistened very much"

TABLE I
CATEGORIES OF EVALUATION EXPRESSION AGAINST "FACE LOTION"

| Rough Categories | Midium Categories | Detailed Categories |
|---|---|---|
| Cost performance | Cost performance | Cost |
| Moisturizing/Penetration | Moisturizing | Keeping Moisture, Moist, Water, Dry/Dry Skin |
| | | Moisturizing, Fresh and young |
| | Penetration | Skin familiarity, Penetration, Suction |
| | Tenseness and elasticity | Elasticity, Springy, Stick to |
| | Tightening the skin | Tightening the skin |
| Whitening care / UV | Whitening care | Whitening care, Dullness, Transparence |
| | UV care | UV care |
| exfoliation & Pores care / Cleansing effect | exfoliation care | exfoliation |
| | Pores care | Pores |
| | Cleansing effect | Cleansing |
| Refreshing feeling / Preventing sebum shine | Refreshing feeling | Refreshing condition, Refreshing feeling |
| | Preventing sebum shine | Tacky, Oil, Shine |
| Refreshing ↔ Thickening | Refreshing ↔ Thickening | Refreshing Texture, Thickening, Sense of use |
| Hypoallergenic | Hypoallergenic | Sensitive skin, Stimulation |
| | Organic | Organic |
| Preventing rough skin | Preventing rough skin | Skin roughness, Trouble |
| | Acne care | Acne care |
| Aging care | Anti-aging | Anti-aging, Beauty ingredient |
| Fragrance | Fragrance | Fragrance, Healing |



Fig. 5.   Constructing the Co-occurrence Keyword-based Dictionary



Fig. 6.   Differences of detecting evaluation expression between Phrase expression-base and Co-occurrence keyword-base

related to an aspect of "Moisturizing", so that the method give a score "7" as an average of their scores based on the phrase expression-based dictionary. Moreover, it includes phrases as "skin irritation issues" related to an aspect of "Hypoallergenic". Thus, the method can give a score "2" based on the dictionary.

Next, we examine a scoring test for non-scored review data based on the constructed dictionary, and compare the result with ground truth data in order to evaluate the effectiveness of the phrase expression-based dictionary. The ground truth data is provided based on not the dictionary but manual detection. We use 16 non-scored reviews and compare numbers of evaluation expressions that are scored by each method in this test.

As a result of the scoring test, a number of evaluation expressions detected manually is 101, whereas a number of evaluation expressions scored based on the phrase expression-based dictionary is 5. That is about only 5% of ground truth data. The reason of such result seems that it is very difficult to construct the phrase expression-based dictionary that can cover various evaluation expression phrases in a large amount of reviews. Because there are many kinds of phrasal expressions even if they are almost same meaning. Therefore, we think that it is necessary to construct not a phrase expression-based dictionary but another dictionary which can cover more evaluation expressions.

*2) Co-occurrence keyword-based Dictionary:* As mentioned in previous section, it is diffifcult for phrase expression-based dictionary to cover most evaluation expression in a lot of reviews. Thus, we try to construct another dictionary using co-occurrence keyword-based evaluation
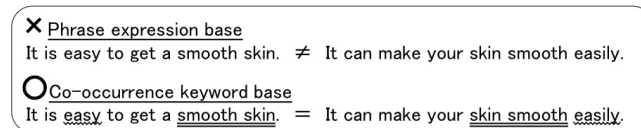
expressions in order to cover wider scope of evaluation expressions.

Fig.5 describes conceptual diagram of constructing the co-occurrence keyword-based dictionary. The procedure of constructing the dictionary is as follows:

1) Analyze phrasal evaluation expressions extracted from reviews.
2) Divide the phrasal expressions into aspect keywords, feature words and degree words.
3) Construct the dictionary by assembling their co-occurrence relations and the evaluation scores.

Fig.6 shows differences of detecting evaluation expression between phrase expression-base dictionary and co-occurrence keyword-base dictionary. As shown in Fig.6, a phrase "It is easy to get a smooth skin" and another phrase "It can make your skin smooth easily" are semantically nearly identical but are different as a phrase. Hence, it may be possible to detect more evaluation expressions based on the co-occurrence keyword-base dictionary than based on the phrase-based dictionary.

*C. Automatic Scoring based on Evaluation Expression Dictionary*

The procedure of automatic scoring against non-scored reviews based on the evaluation expression dictionary is shown below (see Fig.7):

1) System examines a morphological analysis against non-scored review data, and investigates existence or non-existence of aspect keywords as evaluation expression for cosmetic items in the review data.
2) If an aspect keyword exist in it, system checks presence or absence of co-occurrence feature words and
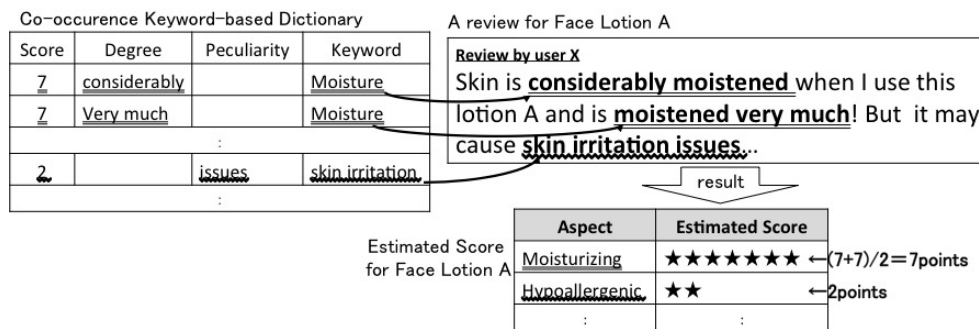
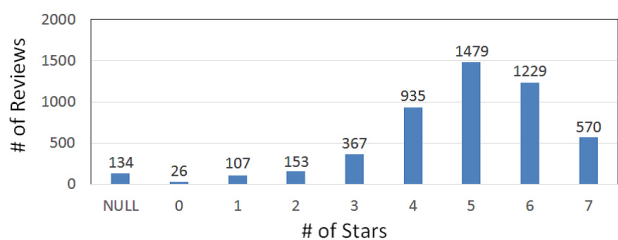Fig. 7.   Automatic Scoring based on the Co-occurrence Keyword-based Dictionary

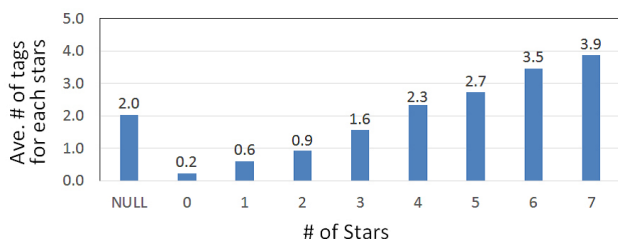

Fig. 8.   Number of Reviews for each Star



Fig. 9.   Average Number of Tags for each Star

degree words co-occurring with the aspect keyword.

3) System make an inquiry on the dictionary to get the score of the evaluation expression based on an aspect keyword, a feature word and a degree word.

4) System achieve "automatic scoring of various aspects of review texts" by aggregating such scores for each aspect in a review.

## IV. EXPERIMENTAL EVALUATION OF AUTOMATIC SCORING USING REAL REVIEW DATA

We examine an experimental evaluation of the automatic scoring method using real review data in order to verify the effectivity of our proposed method. As a first step, we analyze 5,000 reviews randomly extracted from review data for "face lotion" posted at @cosme, for understanding characteristics of the data.

Fig.8 shows a number of reviews for each star (score), and Fig.9 describes an average number of tags for each star (score) in the 5,000 reviews. By the way, reviewers can describe stars from 0 to 7 as score and tags as the effectiveness against cosmetics items at @cosme website.

According to Fig.8, the average number of stars is 4.94 and a distribution of the data looks balanced. According to Fig.9, the average number of tags is 2.75 and we can see that reviewers who give a good score tend to provide more tags.

### A.  Procedure of Experimental Evaluation

In this experiment, we use 10 review data for "face lotion" randomly picked up from 5,000 review as described above, and compare results by following methods:

- Manual scoring method without the Dictionary (as ground truth data).
- Automatic scoring method based on the evaluation expression dictionary (proposed method).

In the case of the manual scoring, evaluators actually read review texts and score them between 0 to 7 stars for 10 aspects of "face lotion" set in advance. The evaluators are 14 people. They are 20's or 30's females.

In the case of automatic scoring, the method scores review texts between 0 to 7 stars for the 10 aspects based on co-occurrence keyword-based dictionary.

The 10 aspects for "face lotion" set for the experiment in advance are shown below:

- Cost performance
- Moisturizing
- Whitening care
- Refreshing feeling/Preventing sebum shine
- Refreshing↔Thickening
- Hypoallergenic
- Preventing rough skin
- Aging care
- Frafrance

### B.  Result of the Experiment

Fig.10 describes results of review scoring against 10 reviews based on co-occurrence keyword-based dictionary and manual scoring by evaluators.

The contents of 10 reviews are different from each other, so that the detected aspects are different. The average score (# of stars) of all aspects by manual scoring is 4.92, and the average score by our proposed method is 4.73. The mean absolute error (MAE) is 0.72.

Scores of the manual scoring tend to a little higher than scores of proposed scoring method. However, the range of the score is from 0 to 7 and MAE is 0.72, so that we may say that the results of automatic scoring by proposed method are quite close to the results of manual scoring as ground truth data. Total number of detected aspects are 48 aspects by manual scoring (ground truth) and 39 aspects by automatic scoring (proposed method). Therefore, the achievement rate of our proposed method against manual scoring is about 81%. The achievement rate of phrase-based dictionary is about 5%
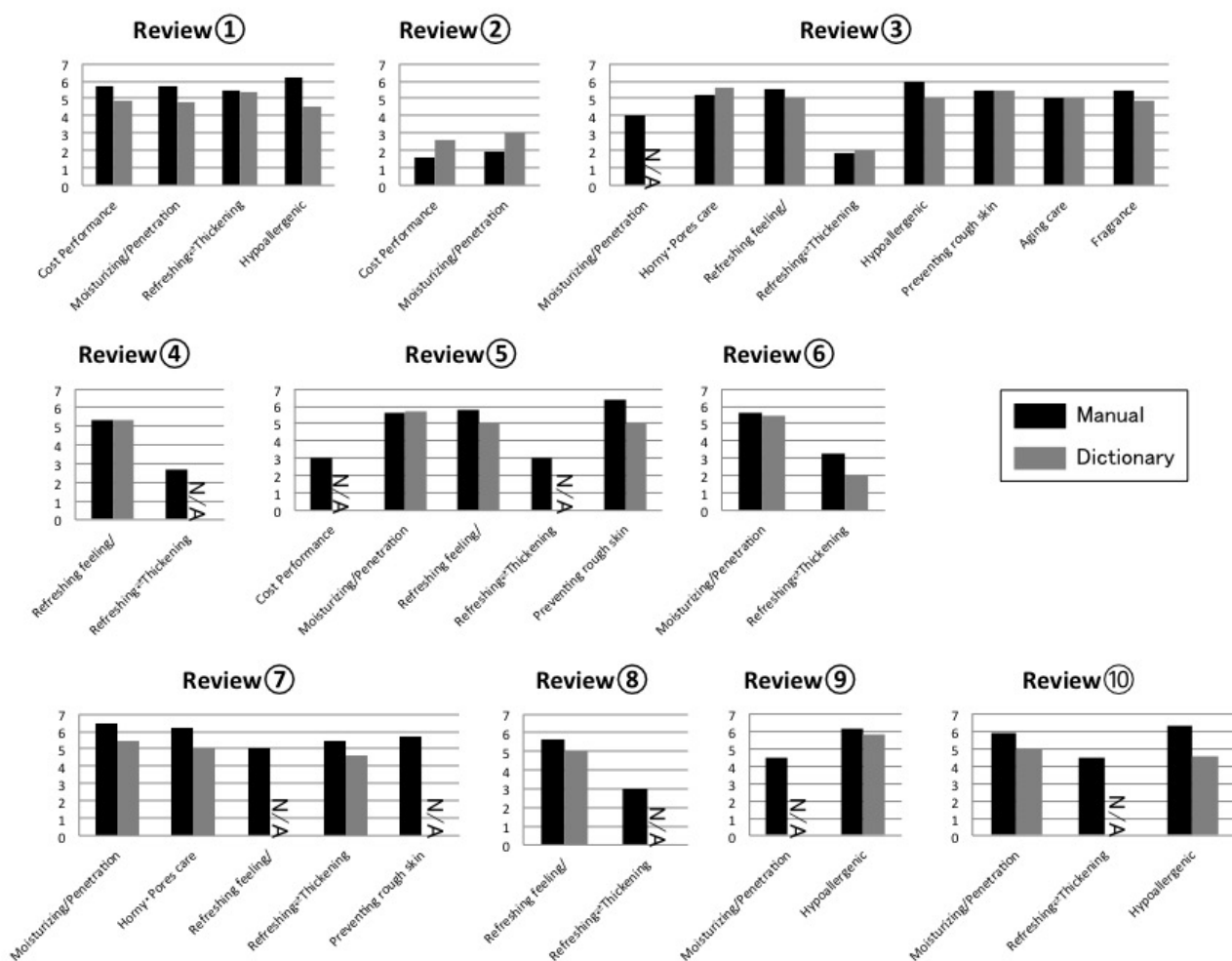
Fig. 10.    Result of Review Scoring based on Co-occurrence Keyword-based Dictionary

as shown in section 3.2.1. Thus, the result of our proposed method based on co-occurrence keyword-based dictionary indicates high potential for detecting aspects of cosmetic items.

There are several "N/A" in Fig.10 by automatic scoring method. However, there are room for improving the result of aspects detection for reviews by updating the dictionary. We will try to analyze larger number of reviews, and then improve and tune the dictionary. Moreover, we will develop a review recommendation system for cosmetic items in our future work.

## V. CONCLUSIONS

In this paper, we presented a method for automatic scoring of various aspects of cosmetic item review texts based on evaluation expression dictionary. In order to realize our proposed method, we constructed two types of evaluation expression dictionaries by extracting and registering evaluation expressions from real review data. Firstly, we constructed a phrase expression-based dictionary. However, it is difficult to cover most evaluation expression in a lot of reviews. Therefore, secondly, we constructed another dictionary using co-occurrence keyword-based evaluation expressions in order to cover the wide scope of evaluation expressions. In order to verify the accuracy of our proposed method, we conducted a simple experiments for the automatic scoring method.

We will improve the dictionary to cover more evaluation expressions, in the future work. A future direction of this study will be to develop a cosmetic item review recommender system which can recommend truly useful reviews for a target user.

## REFERENCES

[1]  @cosme, http://www.cosme.net/, (Accessed 5 January 2016)
[2]  Amazon.com, http://www.amazon.com/, (Accessed 5 January 2016)
[3]  Priceprice.com, http://ph.priceprice.com/, (Accessed 5 January 2016)
[4]  The site data of @cosme (Nov.2015), istyle Inc., http://www.istyle.co.jp/business/uploads/sitedata.pdf (in Japanese), (Accessed 5 January 2016)
[5]  Yumi Hamaoka, Mayumi Ueda and Shinsuke Nakajima, "Extraction of Evaluation Aspects for each Cosmetics Item to Develop the Reputation Portal Site," *IEICE WI2-2012-15*, pp.45-46, 2012.2. (in Japanese)
[6]  Tomohiro Nihongi and Kazuo Sumita, "Analysis and retrieval of the word-of-mouth estimation by structurizing sentences," *Proceeding of the Interaction 2002*, pp.175-176, 2002.3. (in Japanese)
[7]  Ivan Titov and Ryan McDonald, "A Joint Model of Text and Aspect Ratings for Sentiment Summarization," *46th Meeting of Association for Computational Linguistics(ACL-08)*, Columbus, USA, pp.308-316, 2008.
[8]  John O'Donovan, Vesile Evrim, Paddy Nixon and Barry Smyth, "Extracting and Visualizing Trust Relationships from Online Auction Feedback Comments.," *International Joint Conference on Artificial Intelligence (IJCAI'07, Hyderabad, India, January 2007*.
[9]  Byunkyu Kang, Nava Tintarev and John O'Donovan, "Inspection Mechanisms for Community-based Content Discovery in Microblogs" *IntRS'15 Joint Workshop on Interfaces and Human Decision Making for Recommender Systems (http://recex.ist.tugraz.at/intrs2015/) at ACM Recommender Systems 2015. Vienna, Austria. September 2015*