

# **Automatic Extraction of Emotive and Non-emotive Sentence Patterns**

Michal Ptaszynski Fumito Masui **Kitami Institute of Technology** 

Rafal Rzepka and Kenji Araki Hokkaido University

# ABSTRACT

In this research we focus on automatic extraction of patterns from emotive (emotionally loaded) sentences. We assume emotive sentences stand out both lexically and grammatically and verify this assumption experimentally by comparing two sets of such sentences. We use a novel pattern extraction method based on the idea of language combinatorics. Extracted patterns are applied in a text classification task of discriminating between emotive and nonemotive sentences. The method reached balanced F-score of 76% with Precision equal to 64% and Recall 93%.

#### **PROBLEM DEFINITION**

今日はなんて気持ちいい日なんだ!(What a pleasant day today, isn't it?)

# LANGUAGE COMBINATORICS

**SPEC – Sentence Pattern Extraction arChitecture** 

**Sentence patterns = ordered non-repeated combinations of sentence elements.** 

For 
$$1 \le k \le n$$
, there is  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$  all possible *k*-long

$$\sum_{k=1}^{n} \binom{n}{k} = \frac{n!}{1!(n-1)!} + \frac{n!}{2!(n-2)!} + \dots + \frac{n!}{n!(n-n)!} = 2^{n} - 1$$

Normalized pattern weight

$$w_j = \left(\frac{O_{pos}}{O_{pos} + O_{neg}} - 0.5\right) * 2$$



patterns, and

This sentence contains a pattern: なんて \* なんだ! (What a \* isn't it?) 1. This pattern cannot be discovered with n-gram approach.

2. This pattern cannot be discovered if one doesn't know what to look for.

Need to find a way to extract frequent patterns from corpora.

#### DATASET

91 sentences close in meaning, but different emotional load (50 emotive, 41 non-emotive) gathered in an anonymous survey on 30 people of different background (students, businessmen, housewives). Emotive Non amotive

	Emotive	Non-emolve
Examples:	高すぎるからね Takasugiru kara ne	高額なためです。 Kougaku na tame desu.
	'Cause its just too expensive	Due to high cost.
	すごくきれいな海だなあ	きれいな海です
	Sugoku kirei na umi da naa	Kirei na umi desu
	Oh, what a beautiful sea!	This is a beautiful sea
	なんとあの人、結婚するらしいよ	あの日と結婚するらしいです

Nanto ano hito, kekkon suru rashii yo Ano hito kekkon suru rashii desu Have you heard? She's getting married! They say she is gatting married.

**Score for one sentence** 

 $w_j, (1 \ge w_j \ge -1)$ score =

# **EXPERIMENT SETUP**

Preprocessing

Sentence:	今日はなんて気持ちいい日なんだ!				
Transliteration:	Kyōwanantekimochiiihinanda!				
Translation:	What a pleasant day it is today!				
	Preprocessing examples				
1. Tokens:	Kyō wa nante kimochi ii hi nanda !				
2. POS:	N TOP ADV N ADJ N COP EXCL				
3. Tokens+POS:	<i>Kyō</i> [N] <i>wa</i> [TOP] <i>nante</i> [ADV] <i>kimochi</i> [N] <i>ii</i> [ADJ] <i>hi</i> [N] <i>nanda</i> [COP] <i>!</i> [EXCL]				

Pattern List Modification

- All patterns 1.
- 2. **Zero-patterns deleted**
- Ambiguous patterns deleted 3.

All patterns vs. only n-grams

#### Weight Calculation Modifications

- Normalized
- Award length
- Award length and occurrence 3.

Automatic threshold setting

#### **10-fold Cross Validation**

# **EVALUATION EXPERIMENT**

# RESULTS

- <u>Token+POS</u> > Tokenized > POS  $\rightarrow$  algorithm works better on specific elements than more generalized
- <u>Patterns</u> > N-grams (sometimes n-grams get better Precision)
- Length awarded > normalized weight
- Highest results <u>F-score = 0.76</u>, <u>Precision = 0.64</u>, <u>Recall = 0.95</u>
- SPEC slightly worse than ML-Ask [3] (F = 0.79, P = 0.8, R = 0.78)
- SPEC <u>fully automatic</u> > ML-Ask handcrafted



# **DETAILED ANALYSIS**

**Extracted patterns** (Tokenized)

					ine glasses were	
Emotive		Non-emotive		E	Example 2.	
freq.	example pattern	freq.	example pattern		ううん、舞台 <u>が</u> 見え Jun. butai ga mien	
14	、*た	11	い*。	((	Ooh, I cannot see	
12	で	8	し*。	E	Example 3.	
11	ん*。	7	です。	<u></u>	<u> ああ</u> 、おなか <u>が</u> すし	
11	と	6	は*です		<u>\a,</u> onaka <u>ga</u> suita Obb_l'm so bungn	
11	—	6	まし*。			
10	、*た*。	5	ました。		: <b>xample 4</b> . 国なためです	
9	、*よ	5	ます		小山のaku na tame c	
9	、*ん	5	い		ue to high cost.	
8	し	4	です*。	E	Example 5.	
7	ない	3	この*は*	JF J	きれいな海 <u>です</u>	
7		3	は*です。	K	Kirei na umi <u>desu</u>	
6	ん*よ	3	て*ます		his is a beautiful s	
6	、*だ	3	が*た。	E		
6	ちゃ	3	美味しい		_ の本 <u>は</u> とても怖し Cono hon wa totem	
6	よ。	3	た。	T	his book is verv so	
5	だ*。	2	た*、*。	F	xample 7.	
5	に*よ	2	せ		テ日 <u>は</u> 雪が降ってい	
5	が*よ	2	か	k	Kyou <u>wa</u> yuki ga fu	
5	ん	2	さ	lt	is snowing today.	

#### EXAMPLE SENTENCES

#### Example 1.

メガネ、そこにあった<u>ん</u>だ<u>よ</u>。 Megane, soko ni atta <u>n</u> da <u>yo</u>. glasses were over there!)



#### nple 2. <u>ん</u>、舞台<u>が</u>見えない<u>よ</u>。 butai <u>ga</u> mienai <u>yo</u> . I cannot see the stage!)



nple 3. おなか<u>が</u>すいた<u>よ</u>。 onaka <u>ga</u> suita <u>yo</u> . , I'm so hungry)





ï

ï





#### CONCLUSIONS

Presented SPEC - a method for automatic extraction of patterns from emotive sentences. The patterns extracted from a set of emotive and non-emotive sentences were applied in text classification task. Compared different preprocessing techniques (tokenization, POS, token-POS). The best results obtained patterns with both tokens and POS (F-score = 76%, Precision = 64%, Recall 95%). Results for only POS were the lowest. This means the algorithm works better on less abstracted data. The results of SPEC were compared to ML-Ask state-of-the-art affect analysis system. ML-Ask achieved better Precision, but lower Recall. However, as fully automatic, SPEC is more efficient and language independent. Many of the automatically extracted patterns appear in handcrafted databases of ML-Ask, which suggests it could be possible to improve ML-Ask performance by extracting additional patterns with SPEC.

# REFERENCES

- [1] Kaori Sasai. 2006. The Structure of Modern Japanese Exclamatory Sentences: On the Structure of the Nanto-Type Sentence. Studies in the Japanese Language, Vol, 2, No. 1, pp. 16-31.
- [2] Michal Ptaszynski, Rafal Rzepka, Kenji Araki and Yoshio Momouchi. 2011. Language combinatorics: A sentence pattern extraction architecture based on combinatorial explosion. International Journal of Computational Linguistics (IJCL), Vol. 2, Issue 1, pp. 24-36.
- [3] Michal Ptaszynski, Pawel Dybala, Rafal Rzepka and Kenji Araki. 2009. Affecting Corpora: Experiments with Automatic Affect Annotation System -A Case Study of the 2channel Forum-, In Proceedings of The Conference of the Pacific Association for Computational Linguistics (PACLING-09), pp. 223-228.
- [4] C. E. Shannon. 1948. A Mathematical Theory of Communication, The Bell System Technical Journal, Vol. 27, pp. 379-423 (623-656), 1948.
- [5] C. Potts and F. Schwarz. 2008. Exclamatives and heightened emotion: Extracting pragmatic generalizations from large corpora. Ms., UMass Amherst.