# Search for Similar Marks for Detecting Trademark Infringement and Dilution

Shuntaro NAKANO<sup>†</sup>, Hidekazu TANIGAWA<sup>††</sup>, Masaharu MIYAWAKI<sup>†††</sup>, Atsushi YAMADA<sup>††††</sup>, and

Katsumi TANAKA<sup>†</sup>

† Department of Social Informatics, Graduate School of Informatics, Kyoto University Yoshida-honmachi, Sakyo-ku, Kyoto 606–8501 Japan

†† IRD Patent Office, P.O.Box53, OMM Bldg.8F, 1-7-31, Otemae, Chuo-ku, Osaka 540–0008 Japan

††† Ritsumeikan University, College of Law, 56-1 Toji-in Kitamachi, Kita-ku, Kyoto 603-8577 Japan

†††† Advanced Science, Technology & Management Research Institute of Kyoto

134 Chudoji Minamimachi, Shimogyo-ku, Kyoto 600-8813 Japan

E-mail: †{s.nakano,tanaka}@dl.kuis.kyoto-u.ac.jp, ††htanigawa@ird-pat.com, †††mmt23360@law.ritsumei.ac.jp,

††††yamada@astem.or.jp

**Abstract** Trademarks are used to establish brand powers so that products or services are well known in business market. However, there are malicious product/service names (marks) using registered trademark names intentionally or unintentionally and thus these marks cause business damages. This paper proposes a way to detect trademark infringement and dilution by searching for similar marks from the Web. Our proposed method generates and suggests marks that are likely to be considered 'similar' to a given registered trademark. There exist some previous works regarding trademark infringement and dilution, whereas they mainly focused on comparing and judging similarities between multiple registered trademarks. Therefore, they are not so effective to detect suspiciously 'similar' marks in advance. Our proposed method is 1) to generate marks that are likely to be considered are likely to be considered 'similar' in all three types of similarity in the appellation, appearance and concept and 2) to use generated marks as search queries and try to find 'similar' marks from the Web.

**Key words** Trademark Similarity, Examination Guidelines for Trademarks, Similarity in appellation, Similarity in appearance, Similarity in concept

# 1. Introduction

Nowadays, it is indispensable for companies to register names of their products and services as trademarks, making sure that their products and services can clearly be distinguished with other products or services in the market from the consumers' perspectives.

As suggested by Miaoulis and Amato [10], trademarks are also in frequent needs of protection against imitators with the purpose to avoid an ordinary purchaser from being confused, mistaken or deceived regarding the sources of the goods of services. For that reason, trademark registration scheme is available in each country to protect original marks owned by companies or organizations as trademarks so that they are not in danger of imitators. However, in fact companies and organizations are still suffering from intentional or unintentional imitators that are becoming potential risk factors regarding intellectual property management.

Furthermore, Port [13] indicates that, there are also issues of 'dilution' where a third party uses someone's well known trademark without permission (known as 'free riding') and consequently damages the distinctiveness of the trademark that are already obtained.

Based on what is discussed so far, since there are many issues that

can potentially be risk factors for companies with registered trademarks, they should have a clear strategy to protect the trademark distinctiveness from being damages.

To solve this problem, this research aims to assist such companies by offering an information-technology based solution to detect trademark infringement and dilution cases from the Web when one registered trademark is given as the input. In particular, this research focuses on detecting trademark infringement and dilution cases based on trademark-related laws and regulations in Japan.

This research deals with trademark registration system in Japan and thus conforms with 'Examination Guidelines for Trademarks'<sup>1</sup>, published by Japan Patent Office. It defines that when 'judging the similarity of a trademark, decisive elements of the trademark, including its appearance, sound and concept need to be comprehensively taken into consideration'. To follow the principal regarding trademark similarity stated in it, the proposed method in this research handles all three types of trademark similarities; similarity in appellation, similarity in appearance and similarity in concept.

The final goal of the proposed method is to generate a list of

<sup>1</sup> http://www.jpo.go.jp/tetuzuki\_e/t\_tokkyo\_e/tt1302-002.htm

marks that are definitely or potentially similar to the given one registered trademark (hereafter defined as 'suspicious marks') and search the Web using generated marks to discover cases of trademark infringement and dilution.

By utilizing the proposed method of this research, people working at an intellectual property department or equivalent will benefit from being able to quickly find 'suspicious marks' compared to finding them manually, and can proceed to a process of taking legal action against imitators and 'free-riders'. Doing so is necessary to eliminate the potential risks regarding trademarks.

# 2. Related Works

Current researches about trademark similarity can be generally grouped into two groups; one group with researches that focuses on comparing or analyzing multiple trademarks given and other with researches regarding how trademarks are actually received and associated with some impressions in society.

#### 2.1 Comparisons and Analysis of Multiple Trademarks

In the area of trademark comparison and analysis, there are many researches aiming to retrieve trademark images with image features.

In researches of Kim and Kim [19], Wang and Hong [17] and Agrawal et al. [1], they all used Zernike moments of the images as image features and proposed methods to suggest similar figure trademarks when one figure trademark is given as an input.

All of these researches focus on figure trademark similarity and the goals of the researches are to retrieve image trademarks that are potentially similar in appearance by comparing the given one figure trademark and the existing trademark image dataset.

Compared to these previous researches, the proposed method is different in terms of giving one word trademark as an input and handling all three types of similarities in appellation, appearance and concept as 'suspicious marks' generation.

Recently, with the advance of machine learning, some researchers started using machine learning techniques to trademark similarity judgment. For instance, Kawachi and Hiratsuka [7] applied a machine learning technique of neural network to automatically analyze figure trademark similarities in appearance.

There are also researches that analyze several patterns of similarity in appellation. Suga [14] [15] [16] looked at twenty years of trial decision data regarding trademark dispute and analyzed whether increasing or substituting one Japanese mora (unit of sounds that is smaller compared to a syllable) would make the original trademark and a modified mark similar in appellation or not.

Contributions made by Kawachi and Hiratsuka [7] and Suga [14] [15] [16] can be applied to the proposed method of this research to rank a list of marks definitely or potentially similar with one registered word trademark given as the input.

### 2.2 Social Reception and Association of Trademarks

In the previous section, researches comparing and analyzing multiple trademarks are introduced. However, there is one another perspective of researches related in the area of trademarks; that is how trademarks are received in social environment and occasionally associated with one particular impression.

The major approach in this research area is to use the real data representing social reception to measure the power of brand strength. Jamsranjav [6] obtained the review data from users of a cosmetic review website and used the data to measure the power of brand power, defined as consisting of four elements of reception, perception quality, brand image and brand loyalty.

McCarthy [9] is well known as the traditional research about the distinctiveness of the mark and it states that the distinctiveness of a trademark is consisted of its natural distinctiveness and the obtained distinctiveness. As shown in Table 1, the distinctiveness of the trademark is classified into three groups with the type of terms used to name a product or a service, and marketing efforts needed to make the trademark popular in the market.

When arbitrary or fanciful terms are used to name a product or service, it has the strong natural distinctiveness but needs relatively high marketing efforts to make its obtained distinctiveness strong because the coined or fanciful term is not recognized at all when it is named, and vice versa can be said for generic terms.

Table 1: Distinctiveness of marks

INHERENTLY		NON-INHERENTLY	NO
DISTINCTIVE		DISTINCTIVE	DISTINCTIVENESS
No Secondary		Secondary Meaning	No Trademark
Meaning Required		Required	Significance
Arbitrary and Fanciful	Suggestive	Descriptive,Geographic Personal Name	Generic

In this research area, there are also researches focusing on how structural elements of one word trademark are received and associated with one particular impression from average consumers.

Yamade and Haga [18] tried to quantify the subjective degree of similarity of any possible pairs of Japanese '*katakana*' characters with the motivation to eliminate the fatal errors in drug use.

Contributions made by several researches introduced so far can be utilized in the proposed method to construct a model of how consumers feel about one trademark, especially with similarity in appearance and concept.

# 3. Understanding Trademark Similarity

In this chapter, basic concepts about trademarks and trademark similarities used will be clarified for further understandings.

#### 3.1 'Trademark' and 'Mark'

So far, term 'trademark' is used without any clear definition given. However, term 'trademark' actually must be clearly distinguished with similar term of 'mark'. In Japan, regulation regarding trademark is defined in Trademark Act<sup>2</sup>.

<sup>2</sup> http://law.e-gov.go.jp/htmldata/S34/S34HO127.html

In Trademark Act, 'mark' is defined as 'any character(s), figure(s), sign(s) or three-dimensional shape(s), or any combination thereof, or any combination thereof with colors, or any sound(s) and others defined in the government ordinance'.

According to Trademark Act, mark can only be considered as a 'trademark' when it satisfies either of the following two conditions.

• It is used in connection with the goods of a person who produces, certifies or assigns the goods as a business.

• It is used in connection with the services of a person who provides or certifies the services as a business (except those provided for in the preceding condition).

Based on this definition, this research will make clear distinction between 'trademark' and 'mark' for the discussions hereafter.

# 3.2 Types of Trademark

As already mentioned in 3.1, 'mark' is defined as 'any character(s), figure(s), sign(s) or three-dimensional shape(s), or any combination thereof, or any combination thereof with colors, or any sound(s) and others defined in the government ordinance' in Trademark Act of Japan. This definition also applies to when classifying possible type of trademarks. In general, the following list is considered as types of trademark usually used.

- word trademark
- figure trademark
- three-dimensional trademark
- composite trademark (including color)
- sound trademark

In the area of trademark registration, there is one special regulation known as 'standard characters'. 'Standard characters' are set of characters can be used to express word trademarks when an applicant does not wish to express one's trademark with any design elements, such as font style, size or color of the text.

In Japan, the use of 'standard characters' is regulated in Article 5(3) of Trademark Act<sup>3</sup> as 'Where a person desires to register a trademark consisting solely of characters designated by the Commissioner of the Patent Office (hereinafter referred to as "standard characters"), the application shall contain a statement indicating thereof.'.

Furthermore, detailed explanation of how 'standard characters' should be used is available in the corresponding section of 'Examination Guidelines for Trademark'<sup>4</sup>; 'standard characters' that can be used to name word trademarks are also publicly listed as an appendix of 'The Trademark Examination Manual' <sup>5</sup>.

Among several types of trademark generally used, this research focuses on tackling word trademark infringement and 'dilution' problems by giving word trademarks as the input. Also, the proposed method assumes all input word trademarks to be consisted only of 'standard characters'.

#### 3.3 Types of Trademark Similarity

It is generally accepted as given that there are three types of trademark similarity as shown in below.

- Similarity in appellation
- Similarity in appearance
- Similarity in concept

First one is based on an appellation of a trademark, which is about how a trademark is pronounced as the product / service name. Although there are words in many languages including Japanese that can be pronounced in several ways, appellations <sup>6</sup> of a trademark conform with the sound expressions that a owner of that trademark registered with. Therefore, similarity in appellation is measured by comparing registered appellations.

Second one deals with an appearance of a trademark, that is defined as how trademarks look like as shapes. This can be about either of each standard characters' shape in word trademarks or shapes of figure / three-dimentional / composite trademarks. As mentioned in the last section, this research assumes word trademarks consisted of 'standard characters' as the input, so similarity in appearance to be considered in this research will mainly be about shape of each standard characters.

Last similarity with concept, is hard to define compared to similarity in appellation and appearance due to vagueness of concept. Concept of a trademark is an abstract idea that people associate with one particular word or phrase, or even image. This nature of concept makes it difficult to define what similarity in concept is in simple definition. In this research, it only considers concepts associated with one word phrase based on the premise that word trademarks with 'standard characters' are given as the input trademarks.

# 4. Generating and Ranking 'Suspicious Marks'

This chapter will explain the first core element of the proposed method; that is how to generate 'suspicious marks' from the one word trademark given as the initial input for three similarity types and rank by the degree of similarity against the input trademark.

### 4.1 Similarity in Appellation

4.1.1 Generating 'Suspicious Marks'

For similarity in appellation, the detailed regulation is available in Chapter 3, Article 4(1)(xi) of 'Examination Guidelines for Trademarks'<sup>7</sup>. Therefore, the proposed method handling similarity in appellation should strictly conform to the regulation.

However, this regulation is originally for comparing two trademarks and thus can not be used with the original form to generate 'suspicious marks'. In order to make this regulation compatible with generating 'suspicious marks', the proposed method first converts each regulation factors to a rule for 'suspicious marks' generation.

<sup>&</sup>lt;sup>3</sup> http://www.japaneselawtranslation.go.jp/law/detail/?id=45&vm=04&re=01

<sup>&</sup>lt;sup>4</sup> http://www.jpo.go.jp/tetuzuki\_e/t\_tokkyo\_e/pdf/tt1302-002/4.pdf

<sup>&</sup>lt;sup>5</sup> http://www.jpo.go.jp/tetuzuki\_e/t\_tokkyo\_e/pdf/appendix1.pdf

<sup>&</sup>lt;sup>6</sup> A registered trademark sometimes has more than one appellation.

<sup>&</sup>lt;sup>7</sup> http://www.jpo.go.jp/tetuzuki\_e/t\_tokkyo\_e/pdf/tt1302-002/3-10.pdf

#### 4.1.2 Ranking 'Suspicious Marks'

The proposed method uses the co-occurrence probability calculated based on letter n-gram corpora generated from Trademark Gazette CDs with the data offered by Japan Patent Office. The list below shows time periods of CDs used to compile a dataset.

In order to generate letter n-gram corpora, information of trademark itself and its appellations in Japanese 'katakana' are extracted and compiled as the dataset used in the proposed method. The dataset based on Trademark Gazette CDs is consisted of 77,873 registered word trademarks in total.

The letter n-gram corpora are generated using the extracted information of appellation of trademark. There are two n-gram corpora; one as 1-gram corpus and another as 2-gram corpus.

The co-occurrence probability for each generated 'suspicious mark' is calculated using 1-gram corpus and 2-gram corpus. The calculated co-occurrence probability acts as an indicator of how a 'suspicious mark' is likely to be existed as a possible trademark.

### 4.2 Similarity in Appearance

#### 4.2.1 Generating 'Suspicious Marks'

Although 'Examination Guidelines for Trademarks' does not offer concrete examples of similarity in appearance, Amino [2] listed examples available in the judicial precedents and trial decisions.

'Suspicious marks' in terms of 'similarity in appearance' can be generated by replacing a character used when naming products or services with other character that 'looks' similar to the original character. In order to implement character substitutions, the knowledge about which characters are 'similar' in shapes available to use.

#### 4.2.2 Ranking 'Suspicious Marks'

One possible approach to calculate shape similarities between two characters is to compare two characters as images by using some measure and use the calculated value as the parameter indicating shape similarities / differences between two characters.

The proposed method uses Python Imaging Library (PIL) to generate images for three types of characters often used when naming products or services in Japan. Three types of characters are Japanese 'hiragana', 'katakana' and 'kanji'.

In order to calculate similarities between images of characters, HOG(Histogram of Oriented Gradients) feature descriptor by Dalal and Triggs [4] is used to express each of character-based images as a feature vector. HOG feature descriptor break down one input image into several blocks to compute gradient histograms of each blocks and concatenate results into one vector at the end.

Therefore, by computing inner products between HOG feature vectors of two given images using OpenCV, it becomes possible to quantify how two given images are similar; this value is used in the proposed method to eliminate marks that failed to reach the predefined threshold of similarity in appearance from the final output.

#### 4.3 Similarity in Concept

#### 4.3.1 Understanding Similarity in Concept

'Examination Guidelines for Trademarks' states that 'A judgment

on the similarity of a trademark needs to, with consideration given to a class of main users (for example, professionals, senior people, children, women, etc.) of goods or services on which the trademark is used, be made based on attentiveness usually possessed by the user, with consideration given to the state of transaction of the goods or the provision of the services'. Therefore, the differences of attentiveness between classes of main users must also be considered as a part of the proposed method.

In order to conform with this requirement, the proposed method defines 'similarity in concept' by breaking it down into two different large elements with one element also having two child elements, as shown in below and treat them separately and equally when generating and ranking marks that are 'similar in concept'.

- semantic similarity
- dictionary-based similarity
- contextual similarity
- concept building

This classification is based on an idea that two words that are semantically similar are not necessary being 'similar' in terms of 'similarity in concept' regarding trademark. Even though one word or phrase has one or several meanings available, having them does not mean that these meanings are strong enough to be widely accepted as given concepts in everyday real usage.

In other words, there must be some measure to express how a word or phrase is actually received in real society, regardless of how it is defined in dictionary. For this purpose, ideas of 'concept building' are included in a model to define 'similarity in concept'. The parameter of 'concept building' indicates how socially accepted level of two words are similar, and it compares strengths of concepts socially formed for each of two words. This parameter serves important roles when ranking generated 'suspicious marks' in 4.3.3.

For 'semantic similarity', which is a similarity of two words in terms of meanings, it is also further divided into two child elements of 'dictionary-based similarity' and 'contextual similarity'.

'Dictionary-based similarity' is a similarity based on meanings defined in dictionary, and 'contextual similarity' is for meanings actually used in particular context. By combining two of these similarities, the proposed method successfully expresses what it means by 'two words are similar in meanings' into a simplified model.

Furthermore, an issue of 'ambiguity' must be addressed for each of 'semantic similarity' and 'concept building'. The issue of 'ambiguity' is based on a fact that sometimes words can have more than one meanings and it must be detected which meaning is used before comparing meanings or concept strengths of two words.

Based on discussions so far, 'concept similarity' is handled by calculating two parameters of 'semantic similarity' and 'concept building' along with taking ambiguity issues into consideration.

4.3.2 Generating 'Suspicious Marks'

Among aforementioed two elements defining 'similarity in concept', the first element of 'semantic similarity', especially a child element of 'dictionary-based similarity' is considered when generating 'similar' marks and others are considered for ranking.

Usually, word trademarks are made of combining several words and/or coined words. Therefore, by morphologically analyzing the input word trademark, a set of words used in the input word trademark can be obtained. By using synonym databases for each word in the obtained set and replacing one word with corresponding synonyms, marks likely to be 'similar in concept' can be generated.

To do so, Japanese WordNet Synonyms Database<sup>8</sup> distributed as a part of Japanese WordNet Project [5] is used to obtain synonyms for each word. Japanese WordNet Synonyms Database is made by adding Japanese words to original Engligh WordNet [12] and it has information about definition ('gloss') and synonyms of words grouped in a unit called 'synset'. In this dataset, 'synset' corresponds to a meaning of a word unique to that 'synset' ; so one word can have several different 'synset' IDs.

Actual implementation is done by two steps of morphologically analyzing the input word trademark into a set of words using Mecab<sup>9</sup> and substituting each words with all possible synonyms available in Japanese WordNet Synonyms Database. Substitution with possible synonyms is limited to once in one trademark as substituting multiple words, the original form of the input word trademark will gradually be lost and becomes not similar.

4.3.3 Ranking 'Suspicious Marks'

By following processes clarified in last subsection, conceptual 'suspicious marks' for the input word trademark can now be generated. However, some of these generated 'suspicious marks' are not likely be marks similar in concept due to the backgrounds regarding similarity in concept already explained in 4.3.1. Therefore, in order to leave out generated marks that are not in fact similar in concept, some ranking algorithm based on conceptual trademark similarity measure is needed in the proposed method.

To achieve this goal, the proposed method will define a value of conceptual similarity degree  $ConceptSimilarity(w_1, w_2)$  of two words  $w_1, w_2$  as shown in below.

$$ConceptSimilarity(w_1, w_2) = \frac{2}{\frac{1}{SemanticSimilarity(w_1, w_2)} + \frac{1}{ConceptBuilding(w_1, w_2)}}$$
(1)

The above equation is a harmonic mean of *SemanticSimilarity* and *ConceptBuilding* to make sure that two different elements defining 'similarity in concept' are treated with equal weight.

First element of conceptual similarity, *SemanticSimilarity* is a similarity measurement between meanings of two words. To calculate this kind of similarity measurement, some algorithms to express the meaning of a word in a numerical vector is needed so that meanings of two words can be compared.

A major approach to express how words actually mean is

Word2Vec, introduced in the research by Mikolov et al. [11]. Word2Vec uses a text corpus as the input to obtain concurrent words for each word and use to them to express each word as feature vector so that similarity between words can be easily computed. It is known to be effective for expressing the meaning of words as feature vectors with less time needed to learn compared to previous neural network language models. Word2Vec has two models of learning; one is CBOW(Continuous Bag Of Words) and another is Skip-gram.

These two learning models have different learning schemes; while CBOW model is for predicting the current word based on context, Skip-gram model tries to maximize classification of a word based on another word in the same sentence. For this proposed method, an algorithm to express the meaning of one word is needed to use. Therefore, Skip-gram learning model of Word2Vec can be used for the purpose of this research.

However, original Word2Vec model can not handle word ambiguity correctly since it generates one feature vector for one word each. In order to make Word2Vec compatible with handling ambiguous words, Lin et al. [8] proposed a method to distinguish ambiguous words by pre-tagging all words in a given text corpus with types of parts of speeches (grammatical types of words, such as noun, verb, adjective and others) based on the result of morphologically analyzing the input text corpus.

The proposed method in this research also introduced proximity weights in the sum pooling layer of CBOW and this improved algorithm is called as PAS (Proximity-Ambiguity Sensitive) CBOW algorithm. They also scale the proximity weights learned with PAS CBOW to make the sum of them equal to 1. These normalized weights are regarded as a pseudo probability distribution and used for dynamic window size in Skip-gram, as expanded Skip-gram model of PAS Skip-gram model.

However, this ambiguity-aware expansion to Word2Vec learning model is not effective for a problem in this research. That is simply because pre-tagging words with types of speeches often does not offer any distinction between two words with different meanings. For instance, this modified Word2Vec can not make distinctions between 'bank' as 'a land near the river' and 'bank' as 'a financial institution', since both of them are nouns.

Another approach to make Word2Vec compatible of handling ambiguous words correctly is a research by Chen et al. [3]. They approached to these problems by performing word sense disambiguation (WSD) based on different sense vectors created from WordNet.

Although this work seems to be more feasible compared to the one by Lin et al., it still has some problems when applying to the purpose of judging conceptual similarity of Japanese trademarks.

Classifications of synsets and glosses in WordNet are sometimes not so natural as Japanese due to the characteristics of Japanese Wordnet that is compiled by translating original English WordNet. Therefore, it might be expected that sense vector generation and WSD do not work as expected in original English WordNet.

<sup>&</sup>lt;sup>8</sup> http://compling.hss.ntu.edu.sg/wnja/

<sup>&</sup>lt;sup>9</sup> http://taku910.github.io/mecab/

In order to overcome this possible issue, a Word2Vec-based algorithm combining both of some techniques used by Lin et al. and Chen et al. is proposed to calculate *SemanticSimilarity*.

In the proposed method, WSD is processed by mapping each of common nouns in a text corpus to one of WordNet synsets of word, of which cosine similarity between a tf-idf vector of gloss and another tf-idf vector of the line in the text corpus is the largest.

The proposed method uses this ambiguity-aware Word2Vec model to calculate *SemanticSimilarity* defined in the beginning of the current subsection. As Word2Vec model is now aware of ambiguity, it can have several different feature vectors for one word. Therefore, the proposed method considers all possible combinations of feature vectors to calculate similarity between two words. For instance, if word  $w_1$  has m synsets or feature vectors and word  $w_2$ has n synsets or feature vectors, the proposed method will calculate the similarity for all possible m \* n combinations and use the largest similarity value as the similarity between word  $w_1$  and word  $w_2$ .

SemanticSimilarity is a cosine similarity between word vectors of  $w_1$  and  $w_2$ , so it takes the range of -1 to 1. In order to calculate the harmonic mean between SemanticSimilarity and ConceptBuilding at the end, the value of SemanticSimilarity is normalized into range of 0 to 1. This normalized value will be the final value of SemanticSimilarity.

Second element of  $ConceptBuilding(w_1, w_2)$  is based on cases classifications using web hitcounts as shown in Figure 1. Each of  $w_1, w_2, \theta_1$  placed on a horizontal line indicate level of concept popularity in the society. Therefore, in Case 1, concept of word  $w_1$  is considered to having no concept similarity relationship as the popluarity level of  $w_i$  is lower than the threshold value of  $\theta_1$ . If either of  $w_1, w_2$  is smaller than  $\theta_1$ , the value of  $ConceptBuilding(w_1, w_2)$ will be 'undefined'. Another  $\theta$  value of  $\theta_2$ , compared to  $|w_2 - w_1|$ , is the threshold value of concept similarity difference. If the value of  $|w_2 - w_1|$  is equal or larger than  $\theta_2$ ,  $ConceptBuilding(w_1, w_2)$ will be the smallest value of 0.



Figure 1: cases of concept buliding situations

In the actual calculation model of the proposed method, the hitcount value of Bing is used to express the level of concept popularity in the society. However, since the hitcount values tends take a wide range of integer values, the log value with base 2 of the hitcount is used and the inverse of the absolute value  $1/(|log2(w_2) -$   $log2(w_1)|)$  denotes concept similarity relationship.

Finally, as the value of  $SemanticSimilarity(w_1, w_2)$  is normalized to the range of (0,1), same normalization should also be made for  $ConceptBuilding(w_1, w_2)$ . To achieve the normalization, the value of  $ConceptBuilding(w_1, w_2)$  is assumed to have the range of  $range_{ConceptBuilding}$  and using this range,  $ConceptBuilding(w_1, w_2)$  can be normalized into the range of (0,1). This normalized value will be the final value of  $ConceptBuilding(w_1, w_2)$  and the value of  $\theta_1$  is set to 500 and  $\theta_2$  is set to 10,000,000 in the proposed method.

$$range_{ConceptBuilding}$$
(2)  
=  $\left(\frac{1}{|\log 2(\theta_1+1) - \log 2(\theta_1)|}, \frac{1}{|\log 2(\theta_1+\theta_2) - \log 2(\theta_1)|}\right)$ 

The threshold value of concept similarity  $\theta_{concept}$  is determined by calculating *SemanticSimilarity*( $w_1, w_2$ ) for test datasets consisting of trademark pairs that are judged as similar in the count and not similar in the count. These test datasets are compiled using data available by Amino [2] and contain three different types of dataset, hereafter defined as A,B and C respectively.

Dataset A and B are datasets containing 16 pairs of trademarks judged as conceptually similar/not similar by the court respectively.

Usually, for this kind of threshold defining task, preparing two data of similar and not similar is enough by marking each of test data into one of 'true positive', 'false positive', 'false negative' and 'true negative' as shown in Table 2. However, with the task to be solved in this research, the non-similar dataset of B is actually not so well binary-classified. This is due to the nature that mark/trademark similarity cases brought to the court are cases where similarity is disputed between two people because they argue for opposite for similarity and the court is being asked for professional judgement. Therefore, although dataset B does have pairs of trademarks judged as conceptually not similar by the court, it does not necessary mean that pairs in dataset B are similar in concept for sure, in terms of binary-classification. In fact, they are somewhat 'gray' data between 'white (not similar)' and 'black (similar)'.

	Judged Similar	Judged Not Similar
	by the Court	by the Court
Marked Similar	True Positive	False Positive
by	similar by the court	not similar by the court
the Proposed Method	similar by the method	similar by the method
Marked Not Similar	False Negative	True Negative
by	similar by the court	not similar by the court
the Proposed Method	not similar by the method	not similar by the method

Table 2: binary classification of trademark similarity

In order to overcome this issue, another test dataset of C in added to be used alongside with A and B. Test dataset C is made by human changing one of the paired trademark in dataset A to a word that is opposite or not related with the another original trademark in a pair. Dataset C also contains 16 pairs of marks.

By using a unified test dataset containing all of A,B and C, the

value of conceptual similarity threshold  $\theta_{concept}$  is determined. When two marks  $m_1$  and  $m_2$  are given,

• if  $SemanticSimilarity(m_1, m_2)$  is larger than or equal to  $\theta_{concept}, m_1$  and  $m_2$  are similar in concept.

• if  $SemanticSimilarity(m_1, m_2)$  is smaller than  $\theta_{concept}$ ,  $m_1$  and  $m_2$  are not similar in concept.

• if  $SemanticSimilarity(m_1, m_2)$  is undefined, conceptual similarity between  $m_1$  and  $m_2$  is also undefined regardless of the value  $\theta_{concept}$  takes.

Based on this definition and binary classification in Table 2, the value of *Accuracy* is calculated for each of  $\theta_{concept}$  starting at 0.00 and incremented by 0.05 until 1.00 with the equation below. As a result,  $\theta_{concept}$  is defined to be 0.20, which is the smallest value of  $\theta_{concept}$  when *Accuracy* becomes the largest. As trademarks marked 'undefined' are not included in the equation, undefined pairs are separately recorded. In the test dataset, 9 pairs are 'undefined' semantically, and 8 are 'undefined' in concept building.

$$Accuracy_{\theta_{concept}} \tag{3}$$

$$= \frac{tp_{\theta_{concept}} + tn_{\theta_{concept}}}{tp_{\theta_{concept}} + tn_{\theta_{concept}} + fp_{\theta_{concept}} + fn_{\theta_{concept}}}$$

# 5. Search Using Generated 'Suspicious Marks'

#### 5.1 Preparing Search Queries

The proposed method does not use each of generated 'suspicious marks' in the given form; instead a search query is prepared by adding the keyword containing the information about how the input trademark is used with products/services, which users inputed to the proposed system. In an actual application of this research, the search query is prepared in an expression as "generated 'suspicuous mark' AND product/service information keyword' and will be given to the Web search engine of Bing<sup>10</sup>.

#### 5.2 Ranking Searched Web Pages

Using queries prepared in a way explained, searches are performed in Bing to discover webpages containing information related to trademark infringement and dilution cases. Although the proposed method limits number of pages to find for each query to 10 pages, many webpages will be still found. Therefore, some extra efforts to only keep the webpages likely to contain information related to trademark infringement and dilution cases are needed.

### 5.2.1 Used as Goods or Services

The proposed method tries to detect webpages of which generated 'suspicious marks' are used as names of goods or services. If such webpages can be detected, they are more likely to contain information related to trademark infringement and dilution cases compared to webpages of which generated 'suspicious marks' are not used as names of goods or services.

To achieve classification of such webpages, 141 product web

pages in Japanese electronic shopping website of 'Yodobashi Camera<sup>11</sup>' are collected as webpages that are known to contain names of goods or services. These 141 webpages are webpages featured in the top page of 'Yodobashi Camera' as of 16th December, 2016 and there are many product categories of electronics, foods, toys and many others. From the text of these webpages, tags and reserved words of html/css/js/etc. are removed to obtain plain texts.

Using two large text of webpages from 'Yodobashi Camera' and another from newspaper web articles, tf-idf vectors are generated for both of them using Japanese stopwords file distributed as a part of Slothlib<sup>12</sup>. Then, from the tf-idf vector of 'Yodobashi Camera' webpage texts, 300 words with top tf-idf values are extracted as words that represent webpages containing names of goods or services. Words that are below 300 in the ranking of tf-idf values are ignored because 300 seemed to be the border where proper nouns start to appear, as words are examined by human.

### 6. Evaluation

The proposed method of this research will be evaluated separately with two different criteria stated as below.

- (1) Qualities of generated 'suspicious marks'
- (2) Qualities of webpages found by web searches

For evaluating purposes, four word trademarks that each of them is a trademark in trademark pairs that are judged as similar by the court and another four word trademarks that each of them is a trademark in trademark pairs judged as not similar are prepared as total input of eight word trademarks.

#### 6.1 Evaluating Generated 'Suspicious Marks'

Qualities of generated 'suspicious marks' are evaluated by selecting about 0.5% of generated marks for each input word trademark.

Then, two people are asked to judge whether each mark is similar or not similar to the original input word trademark based on their intuition. These two people are people who are chosen intentionally because of their profiles as not being familiar with Trademark Act.

Given T as the input word trademark, S as the marks judged similar by an evaluator, A as the numbers of generated 'suspicious marks' selected for evaluation,  $Precision_T$  of S/A is used as the measurement used to evaluate generated 'suspicious marks'. Table 3 shows the calculated average value of  $Precision_T$  evaluated by two evaluators for eight word trademark given as the input.

### 6.2 Evaluating Web Search Results

In order to evaluate the web pages found by searching with Bing, a registered word trademark of 'メガネの愛眼' was used as an input to the proposed method to generate 'suspicious marks' and search with Bing based using generated marks. Figure 2 shows one example of a webpage found by the proposed method and trademark in logo available in a webpage found. As can be seen in this figure,

<sup>11</sup> http://www.yodobashi.com

<sup>12</sup> http://www.dl.kuis.kyoto-u.ac.jp/slothlib/

<sup>10</sup> https://www.bing.com

Table 3: evaluation of generated 'suspicious marks'

input word trademark $T$	average
レガシィクラブ	0.2757
メガネの愛眼	0.1533
中古車の 110 番	0.1027
ウォークバルーン	0.2578
筑後の寒梅	0.1000
大皿惣菜 ∞ 遊 ∞ 居酒屋	0.1736
和漢研麗姿	0.1005
自然健康館スーパーフコイダン	0.1769
total average of $Precision_T$	0.1675

Figure 2: an example of found webpage and trademark logo

メガネ・サングラスのアイワン 在庫は常時0000年。国内最大級の圧倒え				
CVU <sup>C</sup> OCC -24/32 - 24/34 - 24	お市場サポートダイアル 食 03-6673-7135 2016/2/24実所 法語レビュー放 8,946件 レビュー点脱4.47/5.00 プランド数 200 ほと、在運動は約16,000本1日本一の品紙たに減額11			
MICHAEL KORS Maski Manushira Control Lut Sloung weather Institution Marc 27	URA ASHLEY 後日有 書六作 LACOSTE LAGUNAMOON (2000012) JILLSTUART G U C C I (例2004LDHON 9			
HOTキーワード         超参り1 ジュン・レノンのデッドストックで回         ビジネシーンで           YU磁電く入気1         放気部メオ・東京マランン開始間点         総局量サ           ワングラス         メガネ         SALE         老田県	温ますも1本。 メガネを覚察のマングラス1000 花巻あ州42.5対象メガネ ングラス2000 コスプレイヤー研究後、アンゲーリムモデル *セツリー コンククト レンズ発掘 Q			
当店のコラボモ 数量限定 したいます。 100周年記念 2021年1月 100周年記念 1100周年記念 1100周年記念	モデル第2弾 生産 wed シーン・パレビ マントンには マントン マントンには マントン マントン マントン マントン マントンには マントン マントン マントン マントン マントン マントン マントン マント			
רעטידע- אלגא				

found webpage has a trademark in logo similar to the word trademark of 'メガネの愛眼'.

Since this webpage is difficult or almost impossible to find by searching with 'メガネの愛眼' as the search query, it can be concluded that the proposed method is useful for assisting companies with registered trademarks by offering a solution to detect trademark infringement and dilution from the Web.

### 7. Conclusion and Future Problems

This research proposed a method to generate 'suspicious marks' for three types of trademarks similarities and rank them by the degree of similarity against one given registered word trademark. Furthermore, it also proposed a method to discover trademark infringement and dilution using generated 'suspicious marks' from the Web.

Application of the contribution in this research should not be only limited to discovering a trademark infringement and dilution cases. For instance, there might be a possibility to apply the knowledge of this research for an area of coming up with new product / service name appealing to consumers. Therefore, broader possibilities must be considered when moving this research further forward.

# Acknowledgements

This work was supported by the following project: Grants-in-Aid for Scientific Research (Nos. 24240013) from MEXT of Japan.

#### References

- Deepti Agrawal, Anand Singh Jalal, and Rajesh Tripathi. Trademark image retrieval by integrating shape with texture feature. In *International Conference on Information Systems and Computer Networks*, pp. 30–33. IEEE Comput. Soc, mar 2013.
- [2] Makoto Amino. *Trademark Law*. Yuhikaku Publishing Co., Ltd., sixth edition, 2002.
- [3] Xinxiong Chen, Zhiyuan Liu, and Maosong Sun. A Unified Model for Word Sense Representation and Disambiguation. In *EMNLP2014*, pp. 1025–1035, 2014.
- [4] N. Dalal and B. Triggs. Histograms of Oriented Gradients for Human Detection. In 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), Vol. 1, pp. 886–893. IEEE, 2005.
- [5] Hitoshi. Isahara, Francis. Bond, Kiyotaka. Uchimoto, Masao. Utiyama, and Kyoko. Kanzaki. Development of the Japanese Word-Net. In *the 6th International Conference on Language Resources* and Evaluation, pp. 2420–2423, 2008.
- [6] Baasankhuu Jamsranjav. Visualization of the Brand Power by Wordof-Mouth Data on the WWW (in Japanese). In *Journal of Japan Management Diagnosis Association*, Vol. 7, pp. 348–359, 2007.
- [7] Tomoko Kawachi and Mituyoshi Hiratsuka. Research on Automation of Trademark Similarity Judgment (in Japanese). *IEICE technical report. Social Implications of Technology and Information Ethics*, Vol. 2012-EIP-5, No. 8, pp. 1–7, may 2012.
- [8] Qiu Lin, Cao Yong, Nie Zaiqing, and Rui Yong. Learning Word Representation Considering Proximity and Ambiguity. In Association for the Advancement of Artificial Intelligence, 2014.
- [9] J. Thomas McCarthy. *McCarthy on Trademarks and Unfair Com*petition. Thomson Reuters, fourth edition, 2015.
- [10] G Miaoulis and N D'Amato. Consumer confusion & trademark infringement. *The Journal of Marketing*, Vol. 42, No. 2, pp. 48–55, 1978.
- [11] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. jan 2013.
- [12] George A. Miller. WordNet: A Lexical Database for English. Communications of the ACM, Vol. 38, No. 11, pp. 39–41, 1995.
- [13] Kenneth L. Port. Trademark Dilution in Japan. Northwestern Journal of Technology and Intellectual Property, Vol. 4, No. 2, pp. 228– 254, 2005.
- [14] Fusao Suga. Re-discussing similarity in appellation (1/3) -Chapter 1:influences caused by increasing moras (1/2)- (in Japanese). *Monthly Patent*, Vol. 65, No. 1, pp. 96–106, 2012.
- [15] Fusao Suga. Re-discussing similarity in appellation (2/3) -Chapter 1:influences caused by increasing moras (2/2)- (in Japanese). *Monthly Patent*, Vol. 65, No. 4, pp. 64–70, 2012.
- [16] Fusao Suga. Re-discussing similarity in appellation (3/3) -Chapter 2:inserting and substituting special moras- (in Japanese). *Monthly Patent*, Vol. 65, No. 9, pp. 74–89, 2012.
- [17] Zhenhai Wang and Kicheon Hong. A novel approach for trademark image retrieval by combining global features and local features. *Journal of Computational Information Systems*, Vol. 4, No. February, pp. 1633–1640, 2012.
- [18] Yasuyo Yamade and Shigeru Haga. Subjective evaluation of similarity of appearance of katakana characters in drug names (in Japanese). In *Rikkyo psychological research*, Vol. 50, pp. 79–85. Rikkyo University, mar 2008.
- [19] Yong-Sung Kim and Whoi-Yul Kim. Content-based trademark retrieval system using visually salient features. In *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 307–312. IEEE Comput. Soc, 1997.