

A Regional Food's Features Extraction Algorithm and Its Application

Trung Duc Nguyen

Faculty of Environment and Information Studies

Keio University

5322 Endo Fujisawa Kanagawa 252-8520 JAPAN

*Submitted in partial fulfillment of the requirements
for the degree of Bachelor*

Advisors:

Professor Kiyoki Yasushi

Diep Nguyen-Thi Ngoc

Copyright©2014 Trung Duc Nguyen

Contents

1	Introduction	1
1.1	Background	2
1.2	Challenges and Research Goals	3
1.3	Structure of Thesis	3
2	Food’s Feature-Ingredient Extraction Algorithm	4
2.1	Ingredient Frequency	4
2.2	Ingredient Amount	5
2.3	Ingredient Unique	6
2.4	Featured Index	6
3	Implementation of HUSTLE	8
3.0.1	The Recipe Database	8
3.0.2	Ingredient Frequency	9
3.0.3	Ingredient Amount	9
3.0.4	Ingredient Uniqueness	10
3.0.5	Featured Index	10
4	Application	13
4.1	The System’s Outline	13

4.2	The System's Model	14
4.2.1	The Recipe Suggestion Module Based On Available Materials	14
4.2.2	The Featured Index Calculation Module	15
4.2.3	The Nearest Recipe Selection Module	16
5	Conclusions and Future Works	19
5.1	Conclusions	19
5.2	Future Works	19

List of Tables

3.1	Ingredient Frequency of Ingredients in Kanto region vs Shikoku region	9
3.2	Ingredient Amount of Ingredients in Kanto region vs Shikoku region	10
3.3	Ingredient Amount of Ingredients in Tyubu region vs Kansai region	11
3.4	Ingredient Uniqueness of Ingredients in Japan	11
3.5	Featured Index of Ingredients in Kanto region and Shikoku region	12

List of Figures

4.1	The System's Model with Recipe Suggestion Module, Featured Index Calculation Module and Nearest Recipe Selection Module.	14
-----	--	----

Chapter 1

Introduction

This chapter describes the background of our research and the overview of the problem.

1.1 Background

In this paper, we present a food analysis system to discover the taste of food and understand the featured ingredients in a specific geographical region.

Cooking is the art of making foods. Nowadays, together with the development of technology and the availability of equipment in cooking, many supporting systems are introduced. For example, the cooking support system utilizing built-in cameras and projectors [1], the cooking support system by using ubiquitous sensors [2], the calorie measurement system by image processing [3] or the system which helps inexperienced users in understanding non-professional recipe descriptions [4], etc. However, by using analysis job we can discover the dominant ingredients and tastes in foods and understand how to alter the taste from one to another.

We can observe that in geographical regions that are far apart from each other often have different features and tastes. For example, the Kanto region, which is located in the East of Japan, often has dense taste in its foods, while the foods in Kansai region, which lies in the southern-central region of Japan’s main island Honshu, often has a diluted taste. The reason is because each region has its own special materials for foods and people in these regions have different habits in cooking food. To understand each region’s featured taste we need to answer the following questions: “How can we understand the different features of each region’s food?” and “What effects change the region’s food taste?”. Among the many factors that affect a food’s taste, the combination of materials is a direct and important factor. Each recipe has a list of its own ingredients together with their amount. This leads us to the idea that we could automatically achieve the features of a region’s foods by analyzing the materials. In order to make advantages of recipes in food analysis, we collected recipes in many regions to build a recipe database and propose some analysis algorithms based on

text processing.

1.2 Challenges and Research Goals

We also realize that understanding the region’s featured taste and the preferred materials has an application in supporting cooking activities. For example, imagine someone living in Kanto region who wants to eat some traditional foods in the Kansai region. They know the original recipe but there are some tastes in Kansai region that are not favored. They would prefer that traditional foods with replaced ingredients that are easy for Kanto people to eat. Conversely, someone living in Kanto region might want to try Kanto foods with Kansai taste. Solving this kind of problem means we can build up a system which can help people satisfy their taste. The recipes, which are made by the system, would be flexible and diverse.

In the existing cooking support systems, the methods vary such as image processing, text retrieval, sensing, etc. We use the text processing approach to directly analyze the recipes with their ingredients and amount of ingredients.

1.3 Structure of Thesis

The outline of this paper is as following. The algorithm and the experimental results are introduced in Section 2 and Section 3 respectively. Section 4 describes the web-based application using the proposed algorithm while Section 5 concludes and discusses the remaining problems and future works.

Chapter 2

Food’s Feature-Ingredient Extraction Algorithm

In this section, we propose an algorithm for analyzing the dominant materials which are often used in a region. We define a material in a region to be a featured one if it appears many times with a large amount and be unique among recipes in that region. To evaluate whether it is featured or not, we suppose that the following questions should be answered: “How often, how much, and how unique the material is?”. Respectively, we propose three kind of functions to answer these questions. They have the key role of the metrics for the featured ingredient’s evaluation.

2.1 Ingredient Frequency

The first function named IF (Ingredient Frequency) is used to treat the question “How often does the material appear in a region?”. The higher frequency an ingredient appears in a region, the higher possibility it is the region’s featured ingredient. In each recipe, an ingredient only appears one time. Thus, the time that ingredient appears in the region is the number of the recipes in the region has it as ingredient. Because the database we have from the Internet are often unbalanced, there are some regions that have more recipes than others. Thus to make it independent from the database, we prefer to use the ingredient’s

frequency rather than its appearance times. This function is formed by the number of times the ingredient appears in the region's recipes over the number of total recipes in that region. Let R be the set of all recipes (r) in a region and i be an ingredient which appears in the region. The function is formed as follows:

$$IF(i, R) = \frac{|\{i | i \in r, r \in R\}|}{|R|}$$

Because the IF value is the ingredient's frequency, it takes the value between 0 and 1.

2.2 Ingredient Amount

The ingredient's frequency has little meaning if there is a small amount of it in the recipes. Thus, the taste of a food not only depends on the ingredients, but also the amount of the ingredients. Even when an ingredient has a high value of IF , it might not be the region's featured ingredient. Thus, the second function, IA , is proposed for the question "How much?"

Let r be a recipe in the set of recipes S and ingredient i is in r . We define the mean function $M(i, S)$ be the mean amount of i in S as follows:

$$M(i, S) = \frac{\sum_{i \in r, r \in S} amount(i, r)}{|\{i | i \in r, r \in S\}|}$$

in which $amount(i, r)$ is the amount of ingredient i in recipe r .

We also assume that AR is the set of all recipes in the country regardless of the region it belongs to, while R is the set of all recipes just in a specific region. Thus, $M(i, R)$ calculates the mean amount of ingredient i in the region's recipes (R) while $M(i, AR)$ calculate the mean amount of ingredient i in all the country's recipes (AR). We have the IA function as follows:

$$IA(i, R) = \frac{M(i, R)}{M(i, AR)}$$

Because the IA function calculates the mean of ingredient's amount, it is independent to the frequency of that ingredient. The higher IA value is, the higher possibility it is the region's featured ingredient. Because both numerator and denominator in the formula have the same unit, the IA value is non-unit. Therefore, regardless to the variety of the ingredient's unit, we have a stable metric for evaluating the ingredient's amount.

2.3 Ingredient Unique

The IF and IA functions above might tell us how often an ingredient appears in the region, but this ingredient can often appear in many regions. To be a featured ingredient of a region, the ingredient must satisfy the condition that it appears in the region but doesn't appear in many other regions. We propose the third function IU as follows:

$$IU(i, A) = \log \frac{|A|}{|\{i | i \in a, a \in A\}|}$$

in which i is the ingredient in region a and A is the set of regions.

This function calculates the uniqueness of an ingredient among all the regions. The more often an ingredient appears in different regions the less unique it is. In other words, it is not the featured ingredient of the region. The higher IU value corresponds to higher possibility it is the region's featured ingredient. We use the log scale to make sure the IU values are not too big.

2.4 Featured Index

Featured Index, which is denoted by FI , is the index used to rank ingredients in a region in term of featured ingredient. We realize that these three functions are all proportional to

the rank of the featured ingredient, thus we proposed FI to be the production of these three function's values as follows.

$$FI(i, R) = IF(i, R) \times IA(i, R) \times IU(i, A)$$

The FI function returns the featured index of ingredient i in a region which has a set of recipes R . A is the set of all regions in the country. The ingredients which have the highest FI would be the featured ingredients. On the other hands, the ingredients which have the lowest FI would be considered as the common ingredients for every region.

Chapter 3

Implementation of HUSTLE

This section describes the recipe database and experimental studies on this database by applying the featured ingredient analysis algorithm.

3.0.1 The Recipe Database

We build a recipe database in which recipes are grouped by region. A script written in Python crawls all the recipes from a Japanese cooking website [5]. We chose this website because the recipes are typical foods grouped by regions. The website is only for Japanese recipes, thus we now only have the database for Japanese foods. Each food is characterized by its name, the region it belongs to and its recipe. Each recipe is stored as a map collection in which the ingredient is the key and the couple of amount and unit is the value. Each of the recipes we get from the website is created for various amounts of people. For example, there are recipes for 4 people but there are also recipes for 3 people. Thus we need to normalize the ingredients' amount in each recipe for one person.

There are about 200 recipes over 7 regions in Japan: Kanto, Hokkaido-Tohoku, Shikoku, Tyubu, Kyusyu-Okinawa, Kansai and Tyugoku. We calculate all the above functions for every recipe in Japan, but we only show the experimental results of Kanto and Shikoku within this paper. We chose these two regions because they lie far apart in different islands of Japan. The experimental results are discussed below.

Table 3.1: Ingredient Frequency of Ingredients in Kanto region vs Shikoku region

Kanto region		Shikoku region	
Ingredient	IF	Ingredient	IF
Soy Sauce (醤油)	1.00	Soy Sauce (醤油)	1.00
Miso (みそ)	0.9	Salt (塩)	1.00
Sugar (砂糖)	0.83	Rice (米)	0.83
Sake (酒)	0.83	Sake (酒)	0.67
Salt (塩)	0.67	Green onion (万能ねぎ)	0.50
...
Dried bonito (かつお節)	0.08	Kelp soup (ダシ昆布)	0.16
Pumpkin (かぼちゃ)	0.08	Deep-fried Tofu (油揚げ)	0.16
Kamaage Shirashi (釜揚げしらす)	0.08	Seared bonito (鰹の敲き)	0.16

3.0.2 Ingredient Frequency

Table 3.1 shows that there are some common ingredients which often appear in both Kanto and Shikoku regions such as Soy Sauce (しょうゆ), Sake (酒), Salt (塩),... This is reasonable because we know that these ingredients are common in Japan. Because they often appear in other regions, the IF function is not enough to evaluate the region's featured ingredients. However, it helps us partially understand the habit in using materials in regions. For example, Green onion (万能ねぎ) often appears in Shikoku but not in Kanto region and Sugar (砂糖) often appears in Kanto but not in Shikoku region. This leads us to the idea that typical Kanto foods are often sweeter than Shikoku foods.

3.0.3 Ingredient Amount

Table 3.2 shows the result of the IA value for Kanto and Shikoku region. We can see that most of the IA values are around 1, which means there is not much difference in the way of using an ingredients' amount between Shikoku region and other regions. However, there

Table 3.2: Ingredient Amount of Ingredients in Kanto region vs Shikoku region

Kanto region		Shikoku region	
Ingredient	IA	Ingredient	IA
White radish (大根)	4.27	Shredded seaweed (刻みのり)	6.00
Tempura flour (天ぷら粉)	3.20	Carrot (にんじん)	3.95
Shredded seaweed (刻みのり)	3.00	Tempura flour (天ぷら粉)	3.20
...
Taro (里芋)	0.02	Sweet potato (さつまいも)	0.06
Cake flour (薄力粉)	0.02	Chicken thigh (鶏もも肉)	0.05
Field mustard (菜の花)	0.02	Sushi vinegar (すし酢)	0.05

are some interesting results. For example, in Kansai region, the mean amount of pepper (こしょう) is 11 times greater than the mean amount of total peper in Japan. See details in Table 3.3.

3.0.4 Ingredient Uniqueness

Table 3.4 reflects the fact that the common ingredients such as Salt (塩), Sweet cooking wine (みりん), Ginger (しょうが), Soy sauce (しょうゆ) appear in almost every regions in Japan while the ingredients such as Peanut (落花生) and Chive (あさつき) are not too common and mostly appear in only one region. The ingredients which have the *IU* value of 0 appear in every region.

3.0.5 Featured Index

The Featured Index (*FI*) is the main metric we use to evaluate the regions' featured ingredients. Table 3.5, which is the experimental result of *FI* calculation for Kanto vs Shikoku region, shows us some interesting information. For example, Natto (納豆) is the ingredient

Table 3.3: Ingredient Amount of Ingredients in Tyubu region vs Kansai region

Tyubu region		Kansai region	
Ingredient	IA	Ingredient	IA
Pork loin (豚ロース肉)	25.50	Pepper (こしょう)	11.00
Seaweed (刻みのり)	6.00	Sweet cooking wine (みりん)	6.88
Green onion (長ねぎ)	3.72	Soy sauce (醤油)	5.49
Onion (玉ねぎ)	3.60	Green onion (長ねぎ)	3.72
...
Taro (里芋)	0.03	Milk (牛乳)	0.04
Cake flour (薄力粉)	0.02	Minced chicken (鶏ひき肉)	0.03
Pepper (こしょう)	0.01	Soup (だし汁)	0.01

Table 3.4: Ingredient Uniqueness of Ingredients in Japan

Ingredient	IU
Peanut (落花生)	2.80
Chive (あさつき)	2.80
...	...
Salt (塩)	0.00
Sweet cooking wine (みりん)	0.00
Ginger (しょうが*)	0.00
Soy sauce (しょうゆ)	0.00

Table 3.5: Featured Index of Ingredients in Kanto region and Shikoku region

Kanto region		Shikoku region	
Ingredient	FI	Ingredient	FI
Natto (納豆)	0.60	Kelp (昆布)	1.40
Dried radish (切干大根)	0.47	Sea bream (鯛の切り身)	0.94
Saury (さんま)	0.47	Ponzu suace (ポン酢)	0.90
...
Vineger (酢)	0.00	Sweet cooking wine (みりん)	0.00
Shredded seaweed (刻みのり)	0.00	Egg (卵)	0.00
Wine (酒)	0.00	Wine (酒)	0.00
Ginger (しょうが)	0.00	Rice (米)	0.00

which has the highest FI value in Kanto region. This means Natto (納豆) is possibly the featured ingredient of Kanto region. In Shikoku region, Ponzu sauce (ポン酢) is also often used for Shikoku's foods. The FI of the same ingredient for different regions might differentiate but we figure that if an ingredient ranks high in one region, it cannot rank high in any other regions. The same thing is true for the low-rank ingredients.

Chapter 4

Application

Using the algorithm we propose a system that will help cooking people transform the typical region's food from the original recipe to a new one that has a typical taste of specified region. For convenience and wider use, we develop this system as a web-based system. The system's outline and the model are described below.

4.1 The System's Outline

The system has two main functions:

- Suggesting possible recipes from the set of available materials inputed by the user.

When people cook, they might already have many materials available in their house such as pepper, chili, chicken, etc. But they have no idea which food is the best choice to cook. Thus, we provide a system which has an extra function that accepts available materials inputed by users and then searches in the recipe database for recipes that are suitable for the inputed materials. "Suitable" means the number of extra-buy materials are the least. The suitable recipes will be shown in order; the smaller the number of extra-buy materials there are, the higher rank that recipe will be.

- Transforming a recipe so that it has a specific region's taste.

This is the most important function of the system. It uses the algorithm to extract the featured materials of the specified region and then transforms the original recipe to the new one.

Based on these two functions we divide the system into three modules. These three modules are shown in the middle of Fig. 4.1, represented by three rectangle boxes. The model of the system is described in the next subsection.

Figure 4.1: The System's Model with Recipe Suggestion Module, Featured Index Calculation Module and Nearest Recipe Selection Module.

4.2 The System's Model

4.2.1 The Recipe Suggestion Module Based On Available Materials

This module responds to the first function of the system, suggesting the possible recipes based on available materials inputed by the user. The input of this module is a set of available materials that the user has. It accesses the recipes database during the calculation and its output will be the list of the recipes which include most of the inputed materials.

This output is passed to the Featured Index Calculation Module as shown in Fig. 4.1. After inputting available materials, the system will search in the recipe database for the most suitable recipes and show them in rank order. The pseudocode is shown as below.

```

for  $recipe \in recipes$  do
     $recipe.lack \leftarrow |recipe| - |recipe \cup inputted\ materials|$ 
end for

    sort the recipes by recipe.lack

    return  $recipes$ 

```

4.2.2 The Featured Index Calculation Module

The user selects one of the recipes recommended by the Recipe Suggestion Module. Then selects the region which they want to transform the recipe in order to have that region's taste. This module applies the region's Featured Materials Extracting Algorithm and outputs the list of Featured Index for all materials in the region then stores them in the *FI* Meta Database as shown in Fig. 4.1. Because we are not using all of the lists to extract the featured materials, we only look at two kinds of the following materials:

- The top rank *FI* materials.

These materials are the materials which are often used in the desired region, but not in other regions.

- The bottom rank *FI* materials.

These materials are the most common materials which are used in almost all regions, but with different amounts.

We use both kinds and combine them with the materials appearing in the original recipe. The result is the list of materials and their amount for the food. The output should look in

the shape as follows:

- Onion 2 (original)
- Lemon 1/2 (original)
- ...
- Natto 100g (top *FI*, newly added, region's average)
- Sugar 100g (bottom *FI*, newly added, region's average)

Among the bottom rank *FI* materials, we only take the materials which are already in original recipe to apply into a new recipe. Among the newly added materials we use the average amount of them in the region. The output of this module is passed to the Nearest Recipe Selection Module.

4.2.3 The Nearest Recipe Selection Module

The output of Featured Index Calculation Module gives us the list of materials and their amount which is suitable for the region's taste. But it doesn't mean that we could use that list to make food. If we immediately apply the list of ingredients with the associated amount, we may have a wrong solution. This is because the newly added ingredients and their associated amounts are just the mean value of ingredients in the region. In result, there is the possibility of a bad tasting food. Thus, we propose to search in the region the nearest recipe in term of ingredients and amount. Then apply the suitable ingredients and its amount in that recipe to our food.

Consider the list of materials as a vector. We calculate the similarity between the region's recipe and the average output above. Because we currently have ingredients and their amounts, there is the problem that the unit of ingredient's amounts are different and we

cannot calculate the similarity. Thus we need to normalize these units. The alternative, we propose, is taking the fraction between the recipe's amount and the average amount all over the country. This gives us the values that are unit-independent, therefore usable for the similarity calculation. There are various methods to calculate the similarity between two vectors [6–8]. Among of these methods, Cosine similarity and Euclidean distance are the most famous methods. In this paper, we use the Euclidean distance, therefore the minimum value is adapted. The details of the algorithm is shown below. $X(x_1, x_2, \dots, x_m)$ with $m \in N$ is the list outputed by the FI Calculation Module and $Y(y_1, y_2, \dots, y_n)$ with $n \in N$ represents a list in the lists of the specified region's recipes

```

for ingredient  $\in$  recipe  $X$  do
     $x_i \leftarrow \frac{\text{amount}}{\text{average amount in the country}}$ 
end for

 $min \leftarrow \infty$ 

for recipe  $\in$  region's recipes do
    for ingredient  $\in$  recipe  $Y$  do
         $y_i \leftarrow \frac{\text{amount}}{\text{average amount in the country}}$ 
    end for
     $similarity \leftarrow \sqrt{\sum_{i=k}^l (x_i - y_i)^2}$ 
    if  $similarity < min$  then
         $min \leftarrow similarity$ 
    end if
end for

return recipes

```

Note that though X and Y don't have to have the same dimensions but we only select the

ingredients i which appears both in X and Y to calculate the similarity. x_i and y_i in which $i \in [k, l]$, are the amounts of ingredient i in X and Y respectively.

Chapter 5

Conclusions and Future Works

5.1 Conclusions

In this paper, we have presented the regional foods' features extracting algorithm and the experimental results. The experimental results partially reflect the featured ingredients in regions.

To show the feasibility and applicability of our algorithm, we build the cooking support system that helps cooking people transform the original recipes to have a featured taste of another regions.

In this paper, the recipes from 7 regions of Japan are used. However, the proposed algorithm is scalable to adapt to any recipe from any region in the world. As a future work, we intent to develop a multilingual translation recipe with the function that can transform recipes between countries.

5.2 Future Works

In fact, the seasoning and non-seasoning ingredients affect the food's taste in different ways. Thus, we also develop the research in that direction, analyze food's features by applying different methods for seasoning ingredients and non-seasoning ingredients.

Acknowledgments

Foremost, I would like to express my great appreciate to my supervisor, Professor Hideyuki Tokuda, Professor Hideyuki Tokuda who gave me advice, guidance and encouragement from the beginning when I has just joined to Hide.Tokuda Laboratory.

I would like to thank Professor Jun Murai, Osamu Nakamura and Kenji Takeda, Associate Professor Keisuke Uehara, Hiroyuki Kusumoto, Jin Mitsugi, Rodney D. Van Meter, and Kazuki Takashio, Assistant Professors Massaki Sato, Noriyuki Shigechika and Jin Nakazawa as well.

I am extremely thankful for Dr.Takuro Yonezawa, who has always supported and guided me when I has just joined to CPSF research group of Hide.Tokuda Laboratory. He also helped me analyze the meaning of my research, clarify my research goal and guided me to write technique papers.

I would like to thank my advisor, Mr.Takuya Takimoto, who has always taken care of me since I joined Hide.Tokuda Laboratory. He gave me a lot of advice and support when I am not sure about my research motivations. I would like to thank Mr.Tomotaka Ito and Mr.Ogawa Masaki for a listening and giving me many useful discussion points and advices. I would like to thank members of CPSF research group in Hide.Tokuda Laboratory for their great friendship and support.

I would like to thank Mr.Tomotaka Ito, Mr.Ogawa Masaki, Ms.Vu Le Thao Chi, Ms.Nguyen Thi Ngoc Diep, my friend Ms.Vu Thi Thai Ha, Ms.Luu Thanh Huong, for their comments

for the thesis structure and my English writings.

I send my special thank to my experiment supporters: Takuya Takimoto, Yuuki Nishiyama, Teruaki Ishiguro, Hiroki Shoji, Yutaro Kyono, Nguyen Anh Tien, Do Trung Kien, Nguyen Tien Thanh, Tran Duc Thang Nguyen Thanh Tung, Nguyen Trung Duc, Tran Ngoc Anh and Dinh Hoang Long.

And last but not least, I am heartily grateful to my family. This is the first time I have left my family to start a new life in Japan, many things happened, success and failure, sadness and happiness, tears and smiles. But they still and always are there, giving me strength to follow all goals of my life.

January 4, 2014

Trung Duc Nguyen

Publications

Long Talk

Trung Duc Nguyen, Diep Thi-Ngoc Nguyen, Yasushi Kiyoki. A Regional Food's Features Extraction Algorithm and Its Application. Workshop on Cooking and Eating Activities in conjunction with ACM Conference on Multimedia. Oct 21, Barcelona, Spain.

Poster

SFC OPEN RESEARCH FORUM 2014.

Bibliography

- [1] S. Morioka and H. Ueda, “Cooking support system utilizing built-in cameras and projectors,” in *MVA2011 IAPR Conference on Machine Vision Applications*, pp. 84–89, TeX Users Group, June 13-15 2011.
- [2] Y. Nakauchi, T. Suzuki, A. Tokumasu, and S. Murakami, “Cooking procedure recognition and support system by intelligent environments,” in *RIISS '09*, TeX Users Group, 2009.
- [3] G. Villalobos, R. Almaghrabi, P. Pouladzadeh, and S. Shirmohammadi, “An image processing approach for calorie intake measurement,” 2012.
- [4] I. Ide, Y. Shidochi, Y. Nakamura, D. Deguchi, T. Takahashi, and H. Murase, “Multimedia supplementation to a cooking recipe text for facilitating its understanding to inexperienced users,” in *IEEE International Symposium on Multimedia*, TeX Users Group, 2010.
- [5] Wiki, “Typical foods by area.” <http://www.s-recipe.com/>.
- [6] Wiki, “Cosine similarity.” http://en.wikipedia.org/wiki/Cosine_similarity.
- [7] Wiki, “Euclidean distance.” http://en.wikipedia.org/wiki/Euclidean_distance.
- [8] G. Qian, S. Sural, Y. Gu, and S. Pramanik, “Similarity between euclidean and cosine angle distance for nearest neighbor queries,” in *Proceedings of the 2004 ACM symposium on Applied computing*, SAC '04, (New York, NY, USA), pp. 1232–1237, ACM, 2004.

