

Marcin Hernes  
Krystian Wojtkiewicz  
Edward Szczerbicki (Eds.)

Communications in Computer and Information Science

1287

# Advances in Computational Collective Intelligence

12th International Conference, ICCCI 2020  
Da Nang, Vietnam, November 30 – December 3, 2020  
Proceedings

 Springer



Editorial Board Members

Joaquim Filipe 

*Polytechnic Institute of Setúbal, Setúbal, Portugal*

Ashish Ghosh

*Indian Statistical Institute, Kolkata, India*

Raquel Oliveira Prates 

*Federal University of Minas Gerais (UFMG), Belo Horizonte, Brazil*

Lizhu Zhou

*Tsinghua University, Beijing, China*


More information about this series at <http://www.springer.com/series/7899>

Marcin Hernes · Krystian Wojtkiewicz ·  
Edward Szczerbicki (Eds.)


# Advances in Computational Collective Intelligence

12th International Conference, ICCCI 2020  
Da Nang, Vietnam, November 30 – December 3, 2020  
Proceedings

*Editors*

Marcin Hernes   
Wrocław University of Economics  
and Business  
Wrocław, Poland

Krystian Wojtkiewicz   
Wrocław University of Science  
and Technology  
Wrocław, Poland

Edward Szczerbicki   
University of Newcastle  
Newcastle, Australia

ISSN 1865-0929                      ISSN 1865-0937 (electronic)  
Communications in Computer and Information Science  
ISBN 978-3-030-63118-5              ISBN 978-3-030-63119-2 (eBook)  
<https://doi.org/10.1007/978-3-030-63119-2>

© Springer Nature Switzerland AG 2020

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG  
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

# Contents

## Data Mining and Machine Learning

Rule Induction of Automotive Historic Styles Using Decision Tree Classifier . . . . .	3
<i>Hung-Hsiang Wang and Chih-Ping Chen</i>	
Deep Learning for Multilingual POS Tagging. . . . .	15
<i>Alymzhan Toleu, Gulmira Tolegen, and Rustam Mussabayev</i>	
Study of Machine Learning Techniques on Accident Data . . . . .	25
<i>Zakaria Shams Siam, Rubyat Tasnuva Hasan, Soumik Sarker Anik, Ankit Dev, Sumaia Islam Alita, Mustafizur Rahaman, and Rashedur M. Rahman</i>	
Soil Analysis and Unconfined Compression Test Study Using Data Mining Techniques. . . . .	38
<i>Abdullah Md. Sarwar, Sayeed Md. Shaiban, Suparna Biswas, Arshi Siddiqui Promiti, Tarek Ibne Faysal, Lubaba Bazlul, Md. Sazzad Hossain, and Rashedur M. Rahman</i>	
Self-sorting of Solid Waste Using Machine Learning. . . . .	49
<i>Tyson Chan, Jacky H. Cai, Francis Chen, and Ka C. Chan</i>	
Clustering Algorithms in Mining Fans Operating Mode Identification Problem . . . . .	61
<i>Bartosz Jachnik, Paweł Stefaniak, Natalia Duda, and Paweł Śliwiński</i>	
K-Means Clustering for Features Arrangement in Metagenomic Data Visualization. . . . .	74
<i>Hai Thanh Nguyen, Toan Bao Tran, Huong Hoang Luong, Trung Phuoc Le, Nghi C. Tran, and Quoc-Dinh Truong</i>	
Small Samples of Multidimensional Feature Vectors . . . . .	87
<i>Leon Bobrowski</i>	
Using Fourier Series to Improve the Discrete Grey Model (1, 1). . . . .	99
<i>Van-Thanh Phan, Zbigniew Malara, and Ngoc Thang Nguyen</i>	
Studying on the Accuracy Improvement of GM (1, 1) Model . . . . .	110
<i>Van Dat Nguyen, Van-Thanh Phan, Ngoc Thang Nguyen, Doan Nhan Dao, and Le Thanh Ha</i>	

## Deep Learning and Applications for Industry 4.0

An Evaluation of Image-Based Malware Classification Using Machine Learning . . . . .	125
<i>Tran The Son, Chando Lee, Hoa Le-Minh, Nauman Aslam, Moshin Raza, and Nguyen Quoc Long</i>	
Automatic Container Code Recognition Using MultiDeep Pipeline . . . . .	139
<i>Duy Nguyen, Duc Nguyen, Thong Nguyen, Khoi Ngo, Hung Cao, Thinh Vuong, and Tho Quan</i>	
An Efficient Solution for People Tracking and Profiling from Video Streams Using Low-Power Compute . . . . .	154
<i>Marius Eduard Cojocea and Traian Rebedea</i>	
Simple Pose Network with Skip-Connections for Single Human Pose Estimation . . . . .	166
<i>Van-Thanh Hoang and Kang-Hyun Jo</i>	
Simple Fine-Tuning Attention Modules for Human Pose Estimation . . . . .	175
<i>Tien-Dat Tran, Xuan-Thuy Vo, Moahammad-Ashraf Russo, and Kang-Hyun Jo</i>	
Human Eye Detector with Light-Weight and Efficient Convolutional Neural Network . . . . .	186
<i>Duy-Linh Nguyen, Muhamad Dwisnanto Putro, and Kang-Hyun Jo</i>	
<b>Recommender Systems</b>	
Robust Content-Based Recommendation Distribution System with Gaussian Mixture Model. . . . .	199
<i>Dat Nguyen Van, Van Toan Pham, and Ta Minh Thanh</i>	
Incremental SVD-Based Collaborative Filtering Enhanced with Diversity for Personalized Recommendation. . . . .	212
<i>Minh Quang Pham, Thi Thanh Sang Nguyen, Pham Minh Thu Do, and Adrianna Koziarkiewicz</i>	
Collaborative Filtering Recommendation Based on Statistical Implicative Analysis . . . . .	224
<i>Hiep Xuan Huynh, Nghia Quoc Phan, Nghia Duong-Trung, and Ha Thu Thi Nguyen</i>	

## Computer Vision Techniques

Object Searching on Video Using ORB Descriptor and Support Vector Machine. . . . .	239
<i>Faisal Dharma Adhinata, Agus Harjoko, and Wahyono</i>	
An Improved of Joint Reversible Data Hiding Methods in Encrypted Remote Sensing Satellite Images. . . . .	252
<i>Ali Syahputra Nasution and Gunawan Wibisono</i>	
3D Kinematics of Upper Limb Functional Assessment Using HTC Vive in Unreal Engine 4 . . . . .	264
<i>Kai Liang Lew, Kok Swee Sim, Shing Chiang Tan, and Fazly Salleh Abas</i>	
2D-CNN Based Segmentation of Ischemic Stroke Lesions in MRI Scans . . . .	276
<i>Pir Masoom Shah, Hikmat Khan, Uferah Shafi, Saif ul Islam, Mohsin Raza, Tran The Son, and Hoa Le-Minh</i>	
Melanoma Skin Cancer Classification Using Transfer Learning. . . . .	287
<i>Verosha Pillay, Divyan Hirasen, Serestina Viriri, and Mandlenkosi Gwetu</i>	

## Decision Support and Control Systems

Design a Neural Controller to Control Rescue Quadcopter in Hang Status . . .	301
<i>Nguyen Hoang Mai, Le Quoc Huy, and Tran The Son</i>	
Multidimensional Analysis of SCADA Stream Data for Estimating the Energy Efficiency of Mining Transport. . . . .	314
<i>Paweł Stefaniak, Paweł Śliwiński, Natalia Duda, and Bartosz Jachnik</i>	
A Simple Method of the Haulage Cycles Detection for LHD Machine. . . . .	326
<i>Koperska Wioletta, Skoczylas Artur, and Stefaniak Paweł</i>	
Haul Truck Cycle Identification Using Support Vector Machine and DBSCAN Models . . . . .	338
<i>Dawid Gawelski, Bartosz Jachnik, Paweł Stefaniak, and Artur Skoczylas</i>	

## Intelligent Management Information Systems

Data Quality Management in ERP Systems – Accounting Case. . . . .	353
<i>Marcin Hernes, Andrzej Bytniewski, Karolina Mateńczuk, Artur Rot, Szymon Dziuba, Marcin Fojcik, Tran Luong Nguyet, Paweł Golec, and Agata Kozina</i>	



A Model of Enterprise Analytical Platform for Supply Chain Management. . . . .	363
<i>Paweł Pyda, Helena Dudycz, and Paweł Stefaniak</i>	
Identification the Determinants of Pre-revenue Young Enterprises Value . . . . .	376
<i>Robert Golej</i>	
Blockchain Platform Taxonomy . . . . .	389
<i>Andrew A. Varnavskiy, Ulia M. Gruzina, and Anastasiya O. Buryakova</i>	
Brain Tumor Medical Diagnosis. How to Assess the Quality of Projection Model? . . . . .	402
<i>Paweł Siarka</i>	
Meta-learning Process Analytics for Adaptive Tutoring Systems . . . . .	411
<i>Gracja Niesler and Andrzej Niesler</i>	
<b>Innovations in Intelligent Systems</b>	
Visualization of Structural Dependencies Hidden in a Large Data Set . . . . .	427
<i>Bogumila Hnatkowska</i>	
Internet Advertising Strategy Based on Information Growth in the Zettabyte Era. . . . .	440
<i>Amadeusz Lisiecki and Dariusz Król</i>	
An Approach Using Linked Data for Open Tourist Data Exploration of Thua Thien Hue Province . . . . .	453
<i>Hoang Bao Hung, Hanh Huu Hoang, and Le Vinh Chien</i>	
A Literature Review on Dynamic Pricing - State of Current Research and New Directions. . . . .	465
<i>Karol Stasinski</i>	
<b>Intelligent Modeling and Simulation Approaches for Games and Real World Systems</b>	
Sentiment Analysis by Using Supervised Machine Learning and Deep Learning Approaches. . . . .	481
<i>Saud Naeem, Doina Logofătu, and Fitore Muharemi</i>	
EEG Based Source Localization and Functional Connectivity Analysis . . . . .	492
<i>Soe Myat Thu and Khin Pa Pa Aung</i>	
Fitness Function Design for Neuroevolution in Goal-Finding Game Environments . . . . .	503
<i>K. Vignesh Kumar, R. Sourav, C. Shunmuga Velayutham, and Vidhya Balasubramanian</i>	

An Application of Machine Learning and Image Processing to Automatically Detect Teachers' Gestures . . . . .	516
<i>Josefina Hernández Correa, Danyal Farsani, and Roberto Araya</i>	
The Effect of Teacher Unconscious Behaviors on the Collective Unconscious Behavior of the Classroom . . . . .	529
<i>Roberto Araya and Danyal Farsani</i>	
<b>Experience Enhanced Intelligence to IoT</b>	
Situational Awareness Model of IoV Based on Fuzzy Evaluation and Markov Chain . . . . .	543
<i>Pengfei Zhang, Li Fei, Zuqi Liao, Jiayan Zhang, and Ding Chen</i>	
A Framework for Enhancing Supplier Selection Process by Using SOEKS and Decisional DNA . . . . .	558
<i>Muhammad Bilal Ahmed, Cesar Sanin, and Edward Szczerbicki</i>	
An Efficient Approach for Improving Recursive Joins Based on Three-Way Joins in Spark. . . . .	566
<i>Thanh-Ngoan Trieu, Anh-Cang Phan, and Thuong-Cang Phan</i>	
Lambda Architecture for Anomaly Detection in Online Process Mining Using Autoencoders . . . . .	579
<i>Philippe Krajsic and Bogdan Franczyk</i>	
<b>Data Driven IoT for Smart Society</b>	
Biomedical Text Recognition Using Convolutional Neural Networks: Content Based Deep Learning . . . . .	593
<i>Sisir Joshi, Abeer Alsadoon, S. M. N. Arosha Senanayake, P. W. C. Prasad, Abdul Ghani Naim, and Amr Elchouemi</i>	
Pattern Mining Predictor System for Road Accidents. . . . .	605
<i>Sisir Joshi, Abeer Alsadoon, S. M. N. Arosha Senanayake, P. W. C. Prasad, Shiaw Yin Yong, Amr Elchouemi, and Trung Hung Vo</i>	
Artificial Neural Network Approach to Flood Forecasting in the Vu Gia–Thu Bon Catchment, Vietnam . . . . .	616
<i>Duy Vu Luu, Thi Ngoc Canh Doan, and Ngoc Duong Vo</i>	
Ensuring Comfort Microclimate for Sportsmen in Sport Halls: Comfort Temperature Case Study . . . . .	626
<i>Bakhytzhan Omarov, Bauyrzhan Omarov, Abdinabi Issayev, Almas Anarbayev, Bakhytzhan Akhmetov, Zhandos Yessirkepov, and Yerlan Sabdenbekov</i>	

## Applications of Collective Intelligence

- Smart Solution to Detect Images in Limited Visibility Conditions Based  
Convolutional Neural Networks . . . . . 641  
*Ha Huy Cuong Nguyen, Duc Hien Nguyen, Van Loi Nguyen,  
and Thanh Thuy Nguyen*
- Experience Report on Developing a Crowdsourcing Test Platform  
for Mobile Applications . . . . . 651  
*Nguyen Thanh Binh, Mariem Allagui, Oum-El-Kheir Aktouf,  
Ioannis Parissis, and Le Thi Thanh Binh*
- Vision Based Facial Expression Recognition Using Eigenfaces  
and Multi-SVM Classifier . . . . . 662  
*Hla Myat Maw, Soe Myat Thu, and Myat Thida Mon*
- An Effective Vector Representation of Facebook Fan Pages and Its  
Applications . . . . . 674  
*Viet Hoang Phan, Duy Khanh Ninh, and Chi Khanh Ninh*

## Natural Language Processing

- Wordnet – a Basic Resource for Natural Language Processing:  
The Case of plWordNet . . . . . 689  
*Agnieszka Dziob and Tomasz Naskręć*
- KEFT: Knowledge Extraction and Graph Building from Statistical Data  
Tables . . . . . 701  
*Rabia Azzi, Sylvie Despres, and Gayo Diallo*
- Devising a Cross-Domain Model to Detect Fake Review Comments . . . . . 714  
*Chen-Shan Wei, Ping-Yu Hsu, Chen-Wan Huang, Ming-Shien Cheng,  
and Grandys Frieska Prassida*

## Low Resource Languages Processing

- Towards the Uzbek Language Endings as a Language Resource . . . . . 729  
*Sanatbek Matlatipov, Ualsher Tukeyev, and Mersaid Aripov*
- Inferring the Complete Set of Kazakh Endings as a Language Resource. . . . . 741  
*Ualsher Tukeyev and Aidana Karibayeva*
- A Multi-filter BiLSTM-CNN Architecture for Vietnamese Sentiment  
Analysis. . . . . 752  
*Lac Si Le, Dang Van Thin, Ngan Luu-Thuy Nguyen, and Son Quoc Trinh*

**Computational Collective Intelligence and Natural Language Processing**

Causality in Probabilistic Fuzzy Logic and Alternative Causes as Fuzzy Duals . . . . . 767  
*Serge Robert, Usef Faghihi, Youssef Barkaoui, and Nadia Ghazzali*

Enhance Trend Extraction Results by Refining with Additional Criteria . . . . . 777  
*Ei Thwe Khaing, Myint Myint Thein, and Myint Myint Lwin*

Image Captioning in Vietnamese Language Based on Deep Learning Network . . . . . 789  
*Ha Nguyen Tien, Thanh-Ha Do, and Van-Anh Nguyen*

Textual Clustering: Towards a More Efficient Descriptors of Texts . . . . . 801  
*Ayoub Bokhabrine, Ismaïl Biskri, and Nadia Ghazzali*

Artificial Intelligence in Detecting Suicidal Content on Russian-Language Social Networks . . . . . 811  
*Sergazy Narynov, Kanat Kozhakhmet, Daniyar Mukhtarkhanuly, Aizhan Sambetbayeva, and Batyrkhan Omarov*

**Author Index . . . . . 821**



# Robust Content-Based Recommendation Distribution System with Gaussian Mixture Model

Dat Nguyen Van<sup>1(✉)</sup>, Van Toan Pham<sup>1</sup>, and Ta Minh Thanh<sup>1,2</sup>

<sup>1</sup> Research and Development Department, Sun Asterisk, Ha Noi, Viet Nam  
{nguyen.van.dat, pham.van.toan}@sun-asterisk.com

<sup>2</sup> Le Quy Don Technical University, Ha Noi, Viet Nam  
thanhtm@mta.edu.vn

**Abstract.** Recommendation systems play an very important role in boosting purchasing consumption for many manufacturers by helping consumers find the most appropriate items. Furthermore, there is quite a range of recommendation algorithms that can be efficient; however, a content-based algorithm is always the most popular, powerful, and productive method taken at the begin time of any project. In the negative aspect, somehow content-based algorithm results accuracy is still a concern that correlates to probabilistic similarity. In addition, the similarity calculation method is another crucial that affect the accuracy of content-based recommendation in probabilistic problems. Face with these problems, we propose a new content-based recommendation based on the Gaussian mixture model to improve the accuracy with more sensitive results for probabilistic recommendation problems. Our proposed method experimented in a liquor dataset including six main flavor taste, liquor main taste tags, and some other criteria. The method clusters  $n$  liquor records relied on  $n$  vectors of six dimensions into  $k$  group ( $k < n$ ) before applying a formula to sort the results. Compared our proposed algorithm with two other popular models on the above dataset, the accuracy of the experimental results not only outweighs the comparison to those of two other models but also attain a very speedy response time in real-life applications.

**Keywords:** Recommendation · Content-based · Gaussian-mixture-model (GMM) · Distribution-recommendation

## 1 Introduction

Due to the proliferation of internet, it has brought tremendous chance for people's lives. On the other hand, the myriad and abundance of information on the web has determined a rapidly increasing difficulty in finding what we actually need in a way that can fit the best our requirements [2, 7, 11]. Recommendation

---

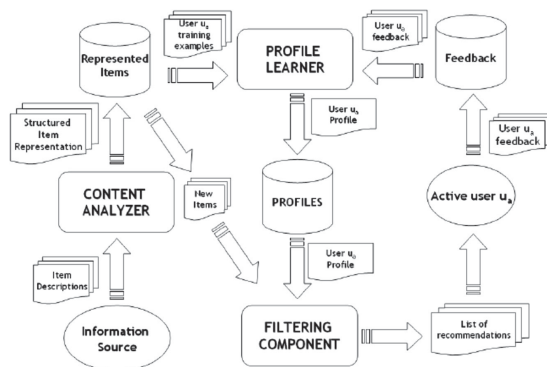
Supported by Sun Asterisk Inc.

© Springer Nature Switzerland AG 2020  
M. Hernes et al. (Eds.): ICCCI 2020, CCIS 1287, pp. 199–211, 2020.  
[https://doi.org/10.1007/978-3-030-63119-2\\_17](https://doi.org/10.1007/978-3-030-63119-2_17)

systems can be effective way to solve such problems without requiring users provide explicit requirements [31,33]. Instead, the system can analysis the content data of item properties, which actively recommend information on users that can satisfy their needs and interests [15,17]. The general content-based architecture is shown in Fig. 1.

Content-based filtering algorithm is widely used because of its simplicity and effectiveness at the begin time of any recommendation systems. According to Pasquale *et al.* [14], there are many benefits reaped from content-based recommendation (CB) systems compared to the other Collaborative Filtering (CF) one such as user independence, transparency, cold-start problems, and so on. Beside, there are still some shortcoming existing as limited content for

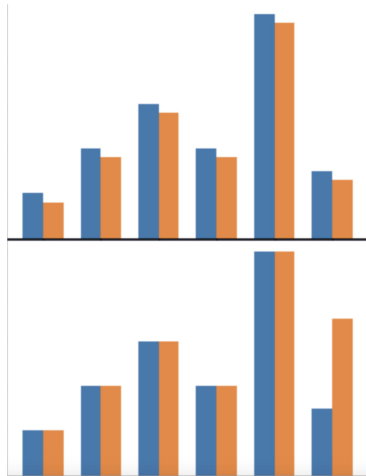
analyzing, over-specialization or lack of rating data of new users and adequate accuracy for some specific problems. Hangyu *et al.* [29] used GMM for CF recommendation algorithm to solve the sparse users rating data. Chen *et al.* [4] proposed a hybrid model, which combines GMM with item-based CF recommendation algorithm and predicted the ratings on items from users to improve the recommendation accuracy. Rui Chen *et al.* [3] took GMM with enhanced matrix factorization to reduce the negative effect of sparse and high dimension data. In the context of music recommender systems, Yoshii *et al.* [30] proposed a hybrid recommender system that combines collaborative filtering using user ratings and content-based features modeled via GMM over MFCCs by utilizing a Bayesian network. However, CF or hybrid systems require behaviour history of users that the reason for the need of CB. Furthermore, CB based on distribution of item features have not been solved yet. A telling of example is using CB for automatically find similar items based on distribution and distance of its features in Fig. 2. These kind of probabilistic problems in recommendation systems is quite different which cannot be solved by usual common methods. Furthermore, the description of content data of items features is sometimes unreliable, inadequate that detrimentally affect to the accuracy of CB systems [11]. Due to two problems mentioned above, we propose a new approach for solving these problems by using GMM [25] to cluster all items into different groups before applying a gaussian filter function (GFF) as a calculation similarity method for sorting results. To demonstrate our effective model, we experiment and compare to two other popular methods, Bag of Word [1] with GFF (BOW + GFF), and GMM with euclidean distance (ED) [13] (GMM + ED). Our propose model not only outperforms the accuracy of the two others, but also get better in prediction time response.



**Fig. 1.** High level architecture of a content-based recommendation system.

The paper is organized as follows. Related work is introduced in Sect. 2 while dataset in Sect. 3. In Sect. 4, the architecture and details of proposed model is given. Experiments and evaluations are shown in Sect. 5. The conclusion are discussed in Sect. 6.

## 2 Related Work



**Fig. 2.** An example between using distribution and distance calculation for distribution recommendation.

We introduce some preliminary knowledge that needs to be used. The following is the detailed information of them.

### 2.1 Content-Based Recommendation

Content-Based Recommendation Systems is one of the most common method in building recommendation systems. The algorithm is born from the idea of using the content descriptions of each item for recommending purposes. It can be divided into two approaches: Analysing the description of item properties only, and building user profile for individuals based on feature's content of items and personal rating data [14,33].

### 2.2 Popular Similarities

In the Content-based algorithm, the similarity calculation method directly affects the accuracy of results. Some similarity calculation methods have been widely used which are listed below:

**euclidean distance:** One of the most popular methods to measure the similarity between two vectors by calculating the sum of square distance of each element respectively in those vectors. Read [13] for more information.

**Cosin:** The main idea is to measure two vectors by calculating the cosine of angle between the two vectors [21].

**Pearson:** The pearson correlation coefficient reflects the degree of linear correlation between two vectors [26],

**Jaccard:** The Jaccard Similarity is often used to compare similarity and different between two finite sample set [19],

### 2.3 Gaussian Mixture Model (GMM)

Gaussian Mixture Model is a function that is comprised of several Gaussians. GMM can t any types of distribution, which is usually used to solve the case

where the data in the same set contains multiple different distributions [5, 24], each identified by  $k \in \{1..K\}$  where  $K$  is the number of clusters of our dataset.

GMM is defined as:

$$p(x) = \sum_{i=1}^k \alpha_i \cdot N(x|\mu_i, \Sigma_i), \quad (1)$$

where  $N(x|\mu_i, \Sigma_i)$  is the  $i^{th}$  component of the hybrid model, which is a probability density function of the  $n$  dimensional random vector  $x$  obeying Gaussian distribution. It can be defined as below:

$$N(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \epsilon^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (2)$$

and

$$\sum_{i=1}^k \alpha_i = 1 \quad (3)$$

We assume that a sample set  $D = \{x_1, x_2, x_3, \dots, x_m\}$  is given that obey gaussian distribution mixture distribution. We use the random variable  $z_j \in \{1, 2, \dots, k\}$  to represent the mixed component of the generated sample  $x_j$ , whose value is unknown. It can be seen that the prior probability  $P(z_j = i)$  of  $z_j$  corresponds to  $\alpha_i (i = 1, 2, 3, \dots, k)$ . According to Bayes' theorem [12], we can get the posterior probability of  $z_j$  as follows:

$$\begin{aligned} p(z_j = i|x_j) &= \frac{P(z_j = i) \cdot p(x_j|z_j = i)}{p(x_j)} \\ &= \frac{\alpha_i \cdot N(x_j|\mu_i, \Sigma_i)}{\sum_{l=1}^k \alpha_l \cdot N(x_j|\mu_l, \Sigma_l)} \end{aligned} \quad (4)$$

In the above formula,  $p(z_j = i|x_j)$  represents the posterior probability of sample  $x_j$  generated by the  $i^{th}$  Gaussian mixture. Assuming  $\gamma_{ij} = \{1, 2, 3, \dots, k\}$  represents  $p(z_j = i|x_j)$ . When the model parameters  $\{(\alpha_i, \mu_i, \Sigma_i) | 1 \leq i \leq k\}$  in the Eq. (4) are known, the GMM clusters divide the sample set  $D$  into  $k$  clusters  $C = \{C_1, C_2, \dots, C_k\}$  [24]. The cluster label  $\lambda_j$  of each sample  $x_j$  can be determined according to equation below:

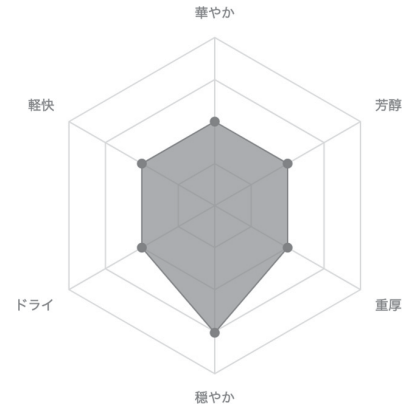
$$\lambda_j = \arg \max_{i \in \{1, 2, 3, \dots, k\}} \gamma_{ji}$$

We get the cluster label  $\lambda_j$  to which  $x_j$  belongs and divide  $x_j$  into cluster  $C_{\lambda_j}$ . The model parameters  $\{(\alpha_i, \mu_i, \Sigma_i) | 1 \leq i \leq k\}$  is solved by applying EM algorithm [16].



### 3 Dataset

Our proposed model is implemented on a dataset about liquor, more specifically, about sake which is one of the most prevalent kind of liquor in Japan. The dataset was collected from Sakenowa<sup>1</sup> being one of the most well-known and reputed website selling the sake<sup>2</sup>. The dataset totally contains 1072 records characterized by 19 properties such as liquor name, liquor brand, year of manufacture, liquor images, liquor flavour tags, liquor six axis flavour taste ( $f_1, f_2, \dots, f_6$  stands for fruity, mel-low, rich, mild, dry and light (Fig. 3). Noticeably, liquor six axis flavour taste and liquor flavour taste would play more important role than the others. The range value of six flavour taste ( $f_1 - f_6$ ) axis is in  $[0, 1]$ , meanwhile the dominant parts belong to  $[0.2, 0.6]$ . The text fields in the dataset all is written behind Japanese form. However, this is a real challenging dataset due to lack of many fields that lead to sparse in data, especially in six main fields  $f_1, \dots, f_6$ . Therefore, our task of recommendation become more difficult and be negatively affect the recommendation results. More specifically, a disappearance or null value of 6-axis fields is greater than 30%, a nearly 2% of null value flavour tags. Further more, many tags value is unreliable, untrust and incorrect that need to be clean and pre-processing (Table 1).



**Fig. 3.** A visualization about 6-axis flavour taste

### 4 Proposed Model

We introduce and explain our proposed model more in detail. As it was mentioned in previous part, we have to return the most similar products based on 19 metadata fields. In particular, 6-axis flavour taste and flavour tags are the main factors mostly affecting to the results both in the sensibility and accuracy side. Therefore, we just select 6-axis flavour taste and flavour tags for better results. The more similar in 6-axis flavour taste, the better results will be.

More detail, we initially use Gaussian mixture model to cluster all items into  $K = \{1, 2, \dots, k\}$  group, then sorting results in each group with each item. Whenever finding top similar items of a item, we just jump up to the group the item belongs to and sort the group's items to return top m similar items. To sort the results, it is also possible to use some popular similarity calculation such as cosine or euclidean distance, but for better accuracy, we use a equation that calculate the distribution weight between two vectors obeying Gaussian

**Table 1.** Dataset blank fields statistic

$f_{1..6}$	Flavour tags	Product name(en)
Float	String	String
30.4%	1.77%	13.4%

<sup>1</sup> <https://sakenowa.com>.

<sup>2</sup> <https://en.wikipedia.org/wiki/Sake>.

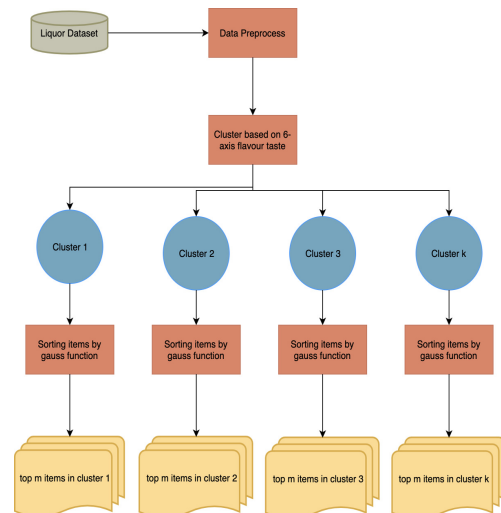
distribution (normal distribution). The results illustrate that the more similar in 6-axis amongst items the bigger weights will be germinated. The flow of our proposed model is shown in Fig. 4.

#### 4.1 Data Pre-processing

It is a fact that text mining is very important in every text-related problems, and CB is not an exception. Previously mentioned, we only choose flavour tags and 6-axis flavour taste as the features for compute the similar between items. The flavour tags are the set of text document written behind Japanese form which require to be cleaned. We convert 6-axis into float and need to do some pre-process techniques for such flavour tags text fields like tokenization, stemmings, stop word removal, find and replace synonyms, lemmatization, and so on [6, 22, 23] before utilizing it. Moreover, the flavour tags field has been splitted into different semantic words, so we disregard the tokenization step and move forward with the other steps.

#### 4.2 Clustering

As we recognize that the final recommendation items depend too much on 6-axis flavour taste and flavour tags. In the common and traditional way, there is a way to build a vector representing for all properties of each item, then utilizing a similarity calculation method like cosine or euclidean to sort and return top  $m$  results. However, in some case, the flavour tags are not enough adequate and precise that adversely affect to the final recommendation. Moreover, there is always an unseen problem of using cosine or euclidean that a compensate between each element of 6-axis flavour taste ( $f_1 - f_6$ ) leads to unequal among those elements ( $f_1 - f_6$ ) of results. Therefore, we decide to group all items based on it's distribution 6-axis flavour taste into different clusters to ensure items which have the same distribution will be in the same cluster that is the foundation for sorting afterwards (Fig. 5).



**Fig. 4.** The model Activities Diagram

### 4.3 Gaussian Function for Sorting

As we have  $K = \{1, 2, \dots, k\}$  clusters, we assume a query item is the center of the cluster we want to find. Our destination is figure out top  $m$  items that have the same distribution as much as possible, so Gaussian filter function (GFF) is the better choice than cosine or euclidean. The Gaussian function equation is defined as follows:

$$G_{kl}(f_{il}, f_{jl}) = \exp - \frac{(f_{il} - f_{jl})^2}{2\sigma_{kl}^2} \quad (5)$$

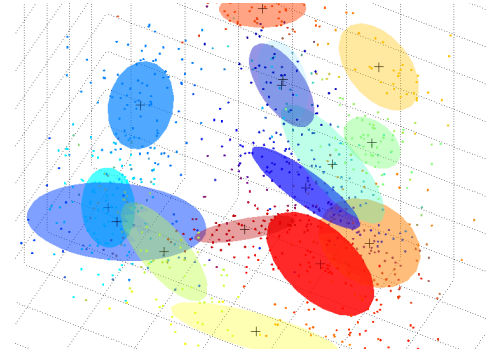


Fig. 5. GMM visualization

where  $G_k(f_{il}, f_{jl})$  is considered as a weight between each pair of element  $l^{th}$  in 6-axis flavour taste of two different items  $(i, j)$  in cluster  $k, l = \{1, 2, \dots, 6\}$ , and  $\sigma_{kl}$  is the standard deviation of the  $l^{th}$  element in 6-axis flavour taste in group  $k$ . Equation for  $\sigma_{kl}$  is defined as below:

$$\sigma_{kl} = \sqrt{\frac{\sum_{i=1}^{n_k} (f_{ilk} - \mu)^2}{n_k - 1}} \quad (6)$$

where  $n_k$  is the number of items belong to cluster  $k$ ,  $f_{ilk}$  is the value of  $l^{th}$  of  $f$  in six flavour taste of  $i^{th}$  item and  $\mu$  is the mean value of all  $f_l$ , in group  $k$ . We calculate  $G(x, y)$  6 times for 6 field  $f_1 - f_6$  for each pair items over all items of a group to sort in descending to find top best results.

### 4.4 Levenshtein Distance for Comparison

The flavour tags also play a quite significant role in the final results. We treat tags as vital as each element in 6 flavour taste. To compare and measure the similarity between two tags of string type, we use a good of levenshtein distance to solve it [8, 32]. The equation of the levenshtein distance is defined below:

$$lev_{a,b}(i, j) = \begin{cases} \max(i, j), & \text{if } \min(i, j) = 0 \\ \min = \begin{cases} lev_{a,b}(i - 1, j) + 1 \\ lev_{a,b}(i, j - 1) + 1, & \text{otherwise} \\ lev_{a,b}(i - 1, j - 1) + 1_{(a_i \neq b_j)} \end{cases} \end{cases} \quad (7)$$

### 4.5 Final Sorting Formula

Combine weight calculating function for 6-axis flavour taste and tags comparison with levenshtein distance (LD), we establish a equation for sorting to get final results as below:

$$S(i, j) = \sum_{k=1}^K \sum_{l=1}^6 G_{kl}(i, j) + lev_{tags}(i, j) \quad (8)$$

where  $G_{kl}$  is the weight function corresponding  $l^{th}$  in 6-axis flavour taste between  $item_i$  and  $item_j$  in cluster  $k$  ( $k = \{1, \dots, K\}$   $K$  groups),  $lev_{tags}(i, j)$  is the levenshtein function to compare tags similarity of those two items. We determine that the bigger  $S(i, j)$ , the better similar between those two items, so we sort by descending order all items of a cluster and return top  $m$  items having bigger  $S(i, j)$  value.

#### 4.6 Proposed Model Pseudo Code

For clearly, we give the proposed model its algorithm execution process to help readers more easily visualize and imagine our entire process. Let see pseudo code below:

---

##### Algorithm 1: Framework model proposal

---

**Input:** number of clusters  $k$

**Output:** Top  $m$  other similar items of each item

**Data:** Dataset  $L$

1. Data pre-processing for text fields
2. Build a matrix for 6-dimension vectors representing for six flavour taste ( $f_1 - f_6$ )
3. Taking the matrix as an input of GMM to train and save corresponding cluster of each item into dataset
4. **for** *item in dataset* **do**
  - Get cluster number of item
  - Find all items that have the same clusters
  - Applying equation  $S(i, j)$  “(8)” for each pairs items
  - Return top  $m$  similar items by sorting decreasingly

**end**

---

## 5 Experiments

To prove the validity of our proposed model, we compare our proposed model to two other popular algorithms widely used in CB systems such as BOW+GFF and GMM+ED. We also illustrate the impact of GMM cluster into the accuracy and the efficiency of Gaussian filter equation in sorting results rather than those of cosine or euclidean distance.

### 5.1 Evaluation Method

The evaluation method of recommendation systems commonly used is Root Mean Square Error (MSE) that is the average of the square errors [27]. It is defined as:

$$MSE = \frac{1}{N} \sum_1^n (r_i - \hat{r}_i)^2, \quad (9)$$

where  $r_i$  is the predicted representing vector item, and  $\hat{r}_i$  is the original representing vector item.

We also use the recommendation results in Sakenowa as the standard measure to compare with our three algorithms because the sake website has so much reputation, popularity, being well-known for commercial purpose in Japan for many years. Beside, the recommendation results of the Sakenowa is also very impressive.

## 5.2 Experimental Analysis

Some experiments will be conducted to verify the impact of GMM, GFF on probabilistic recommendation problem. Our main proposed model was implemented through such steps as data statistic, data cleaning, data missing value filling, clustering all items into different clusters and eventually using Gaussian filter + levenshtein distance to sort the results. To verify the effective impact of GMM and GFF on better prediction, we divided our experiments into three parts. Firstly, we use Bag-of-Word (BOW) [1] algorithm on some properties like flavour tags before applying GFF for sorting results. In the second way, we apply GMM + ED to clarify the influence of GMM. Finally, we implemented our main proposed model to prove the impact of GMM+GFF then give some comparison. All experiments will be unraveled in detail below:

**Experiment 1: BOW+GFF.** The reason for this experiment is to verify the impact on result accuracy of GMM compared to BOW algorithm. Therefore, in the experiment, we will implement BOW algorithm comprised with GFF used for sorting on our liquor dataset. Firstly, we do some data preprocessing for text data like stemming, replace synonyms, filling missing data, etc [22]. As it was mentioned above, all important text fields were written behind Japanese form, so we use some tools offered for Japanese preprocessing like Ginza [9], Janome [10], JapaneseStemmer [18] was inspired by Porter Stemming Algorithm [28], etc. Before using GFF for sorting, we use BOW on these preprocessed properties to find the vector matrix representing for the item. The next step, we feed the vector matrix into  $K$ -nearest neighbors ( $k$ -NN) algorithm using unsupervised  $k$ -NN Scikit-Learn [20] to find top similar items based on these vectors. In these top items, we apply equation  $S$  (8) to get the best similar items.

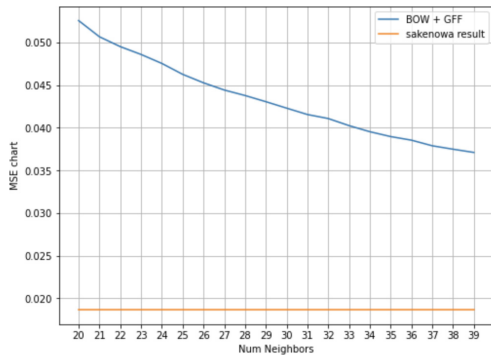
**Experiment 2: GMM+ED.** To demonstrate the impact of GMM, firstly, we still apply some preprocessing steps for text fields as the experiment above. After that, we build a matrix of 6 dimensions representing for 6 flavour taste, then feed it into GMM for training, save all cluster result for each item. Next step, we convert a collection of text flavour tags into a matrix of token counts using CountVectorizer of Scikit-Learn [20] and concatenate along same axis with the

matrix of 6 dimensions for sorting. Finally, to return the best similar items of a given item, we just jump up to the cluster containing it and apply ED for sorting the results and get top best similar items of the given item.

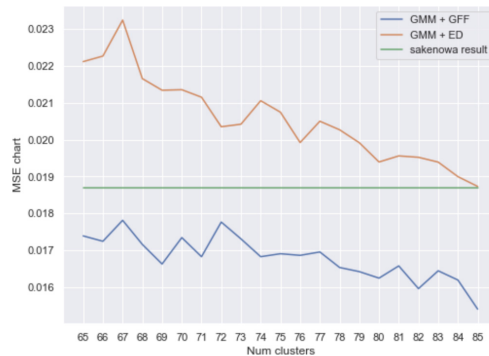
**Experiment 3: GMM+GFF.** Our two above experiment to prove the important role of GMM and GFF in our proposed model. In this experiment, in the first place, we also do preprocess for text fields as same steps in two previous experiments. After that, we build a matrix of 6 dimensions representing for 6 flavour tastes and feed the matrix into GMM for training purpose, then save cluster results for each items. To find top best similar items of a given item, we jump up into the cluster the query item lied in, consider the query item as center then apply (8) equation pair in pair with all items in the cluster, then sorting discerningly to return the top best similar items.

### 5.3 Experimental Results and Comparison

In this section, we compare our proposed algorithm with the results from the Sakenowa website and two other popular CB algorithms. The recommendation results from Sakenowa for each item are returned from an api<sup>3</sup>; therein,  $f_{1...6}$  in the api are the value for each flavour taste, respectively. We conclude that our result accuracy outweighs the Sakenowa and these two algorithm counter parts. Let see some charts below:



**Fig. 6.** MSE applied BOW+GFF

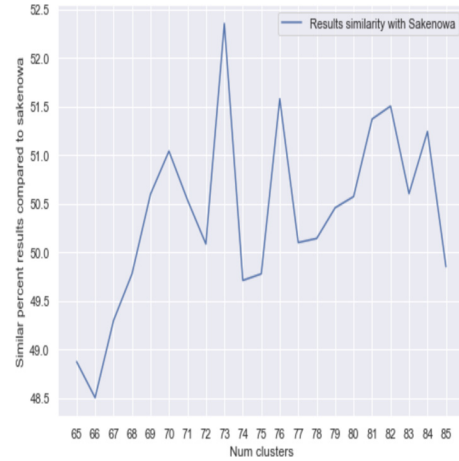


**Fig. 7.** MSE applied GMM+GFF and GMM+ED

All three experiments return top ten(10) best similar items for each item in dataset. In Fig. 6, list values of MSE are shown through an array of number of neighbors ranging from [25–39] in  $k$ -NN algorithm. Despite the tendency of decrease, but it is insignificant and the time response is extremely slow due to bigger number of neighbors.

<sup>3</sup> [https://sakenowa.com/api/v1/brands/flavor?f=0&fv=f\\_1,f\\_2,f\\_3,f\\_4,f\\_5,f\\_6](https://sakenowa.com/api/v1/brands/flavor?f=0&fv=f_1,f_2,f_3,f_4,f_5,f_6).

In Fig.7 the gap of MSE between GMM+ED and GMM+GFF is shown. It is very clearly seen that GMM+GFF generate better results than the other that verify the effect of GMM in sorting results. Both these two experiments show the effect of the number of clusters of GMM ranging from [65–85]. In Fig.8, we compare our prediction results of all items in dataset to the recommendation results from the Sakenowa and construct a list of similarity proportion affected by the number of clusters.



**Fig. 8.** Similar percent statistic compare to Sakenowa

**Table 2.** MSE affected by number of clusters

N clusters	GMM+ED	GMM+GFF	Sakenowa results
65	0.02211	<b>0.01739</b>	0.01868
70	0.02135	<b>0.01734</b>	0.01868
75	0.02074	<b>0.01691</b>	0.01868
80	0.01939	<b>0.01625</b>	0.01868
85	0.01873	<b>0.01541</b>	0.01868

In Table 2 and Table 4, we build a table of statistic of MSE generated from GMM+ED, BOW+GFF, GMM+GFF and recommendation results from Sakenowa. It is matter of fact that our GMM+GFF algorithm outperforms all the others method that demonstrate the effective of our algorithm. Further more, our time response in Table 3 also beat these two others, GMM+ED and BOW+GFF.

**Table 3.** Time response per query

BOW+GFF	GMM+ED	GMM+GFF
0.1856 s	0.0174 s	<b>0.0156 s</b>

**Table 4.** MSE affected by number of neighbors

Num neighbors	BOW+GFF	Sakenowa results
20	0.05254	0.01868
25	0.04624	0.01868
30	0.04228	0.01868
35	0.03895	0.01868
39	0.03709	0.01868

## 6 Conclusion

We have proposed an very effective algorithm for recommendation system using content-based features with GMM. We have applied our proposed method for solving liquor recommendations. Further, our probabilistic-based recommendation systems not only acquire a remarkable prediction accuracy, but also has very speedy prediction time response for real-time application.

## References

1. Bhattacharya, S., Lundia, A.: Movie recommendation system using bag of words and Scikit-learn (2019)
2. Bollen, D., et al.: Understanding choice overload in recommender systems, pp. 63–70 (2010). <https://doi.org/10.1145/1864708.1864724>
3. Chen, R., et al.: A hybrid recommender system for Gaussian mixture model and enhanced social matrix factorization technology based on multiple interests. *Math. Prob. Eng.* 1–22 (2018). <https://doi.org/10.1155/2018/9109647>
4. Fan-sheng, K.: Hybrid Gaussian pLSA model and item based collaborative filtering recommendation. *Comput. Eng. Appl.* (2010)
5. Görür, D., Rasmussen, C.: Dirichlet process Gaussian mixture models: choice of the base distribution. *J. Comput. Sci. Technol.* **25**, 653–664 (2010). <https://doi.org/10.1007/s11390-010-9355-8>
6. Gurusamy, V., Kannan, S.: Preprocessing techniques for text mining (2014)
7. Guy, I., Carmel, D.: Social recommender systems, pp. 283–284 (2011). <https://doi.org/10.1145/1963192.1963312>
8. Haldar, R., Mukhopadhyay, D.: Levenshtein distance technique in dictionary lookup methods: an improved approach. In: *Computing Research Repository - CORR* (2011)
9. Mai, O., Hiroshi, M., Masayuki, A.: Simultaneous learning of ambiguity resolution and dependency labeling for short-unit part of speech usage. In: 25 (2019). [http://www.anlp.jp/proceedings/annual\\_meeting/2019/pdf\\_dir/F2-3.pdf](http://www.anlp.jp/proceedings/annual_meeting/2019/pdf_dir/F2-3.pdf)
10. Janomepy<sub>py</sub> Janome (2019). <https://github.com/mocobeta/janome>
11. Khusro, S., Ali, Z., Ullah, I.: Recommender systems: issues, challenges, and research opportunities, pp. 1179–1189, ISBN 978- 981-10-0556-5 (2016). [https://doi.org/10.1007/978-981-10-0557-2\\_112](https://doi.org/10.1007/978-981-10-0557-2_112)
12. Lee, D.-S., Hull, J., Erol, B.: A Bayesian framework for Gaussian mixture background modeling, vol. 3, pp. III-973 (2003). <https://doi.org/10.1109/ICIP.2003.1247409>
13. Liberti, L., et al.: Euclidean distance geometry and applications. *SIAM Rev.* **56**, 3–69 (2012). <https://doi.org/10.1137/120875909>
14. Lops, P., de Gemmis, M., Semeraro, G.: Content-based recommender systems: state of the art and trends. In: Ricci, F., Rokach, L., Shapira, B., Kantor, P.B. (eds.) *Recommender Systems Handbook*, pp. 73–105. Springer, Boston, MA (2011). [https://doi.org/10.1007/978-0-387-85820-3\\_3](https://doi.org/10.1007/978-0-387-85820-3_3)
15. Linyuan, L., et al.: Recommender systems. *Phys. Rep.* **519**(1), 1–49 (2012). <https://doi.org/10.1016/j.physrep.2012>. ISSN 0370–1573
16. Lu, Y., Bai, X., Wang, F.: Music recommendation system design based on Gaussian mixture model. In: *ICM 2015* (2015)



17. Melville, P., Sindhvani, V.: Recommender systems, pp. 829–838 (2011). [https://doi.org/10.1007/978-0-387-30164-8\\_705](https://doi.org/10.1007/978-0-387-30164-8_705)
18. MrBrickPanda: Japanese Stemmer (2019). <https://github.com/MrBrickPanda/Japanese-stemmer>
19. Niwattanakul, S., et al.: Using of Jaccard coefficient for keywords similarity (2013)
20. Pedregosa, F., et al.: Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* **12**(Oct), 2825–2830 (2011)
21. Philip, S., Shola, P., Abari, O.: Application of content-based approach in research paper recommendation system for a digital library. *Int. J. Adv. Comput. Sci. Appl.* **5** (2014). <https://doi.org/10.14569/IJACSA.2014.051006>
22. Rahutomo, R., et al.: Preprocessing methods and tools in modelling Japanese for text classification, p. 8843796 (2019). <https://doi.org/10.1109/ICIMTech.2019>
23. Rajman, M., Besanc¸on, R.: Text mining: natural language techniques and text mining applications. In: Proceedings of the 7th IFIP Working Conference on Database Semantics (DS-7) (1997). [https://doi.org/10.1007/978-0-387-35300-5\\_3](https://doi.org/10.1007/978-0-387-35300-5_3)
24. Rasmussen, C.: The Infinite Gaussian mixture model, vol. 12, pp. 554–560 (2000)
25. Reynolds, D.: Gaussian mixture models. In: Encyclopedia of Biometrics (2008). [https://doi.org/10.1007/978-0-387-73003-5\\_196](https://doi.org/10.1007/978-0-387-73003-5_196)
26. Sedgwick, P.: Pearson’s correlation coefficient. *BMJ* **345**, e4483–e4483 (2012). <https://doi.org/10.1136/bmj.e4483>
27. Shani, G., Gunawardana, A.: Evaluating recommendation systems, vol. 12, pp. 257–297 (2011). [https://doi.org/10.1007/978-0-387-85820-3\\_8](https://doi.org/10.1007/978-0-387-85820-3_8)
28. Jones, K.S., Willett, P. (eds.): Readings in Information Retrieval. Morgan Kaufmann Publishers Inc., San Francisco (1997). ISBN 1558604545
29. Yan, H., Tang, Y.: Collaborative filtering based on Gaussian mixture model and improved Jaccard similarity. *IEEE Access* **7**, 118690–118701 (2019). <https://doi.org/10.1109/ACCESS.2019.2936630>
30. Yoshii, K., et al.: Hybrid collaborative and content-based music recommendation using probabilistic model with latent user preferences, pp. 296–301 (2006)
31. Zhu, B., Bobadilla, J., Ortega, F.: Reliability quality measures for recommender systems. *Inf. Sci.* **442**, 145–157 (2018)
32. Ziolk, B., et al.: Modified weighted Levenshtein distance in automatic speech recognition (2010)
33. Zisopoulos, H., et al.: Content-based recommendation systems (2008)