

An Explanation Framework for Whole Processes of Data Analysis Applications: Concepts and Use Cases

Hiroshi Ishikawa
Graduate School of Systems Design
Faculty of System Design
Tokyo Metropolitan University
Hino, Tokyo
Email:ishikawa-hiroshi@tmu.ac.jp

Masaharu Hirota
Department of Information Science
Faculty of Informatics
Okayama University of Science
Okayama, Okayama
Email:hirota@mis.ous.ac.jp

Yukio Yamamoto
Japan Aerospace Exploration Agency
Sagamihara, Kanagawa
Email:yamamoto.yukio@jaxa.jp

Masaki Endo
Division of Core Manufacturing
Polytechnic University
Kodaira, Tokyo
Email:endou@uitec.ac.jp

Abstract- The main contribution of the paper is to address the necessity of both macro and micro explanations for Social Big Data (SBD) applications and to propose an explanation framework integrating both, allowing SBD applications to be more widely accepted and used. The framework provides both a macro explanation of the whole procedure and a micro explanation of the constructed model and an explanation of the decisions made by the model. Application systems including Artificial Intelligence (AI) or Data Mining (DM) need reproducibility to ensure their reliability as scientific systems. For that purpose, it is important to illustrate the procedures of the system explicitly and abstractly (that is, macro explanations). This paper has scientific value in that it proposes a data model for that purpose and illustrates the possibility of macro explanations through one use case of social science. Scientists also need to provide evidence that the results obtained by AI or DM are valid. In other words, this paper also has scientific value in that it reveals how the features of the model and concrete grounds for judgment can be explained through two use cases of natural science.

Keywords- social big data; explanayion; data model; data management; data mining.

I. INTRODUCTION

We are surrounded by big data, which are waiting to be analyzed and used. Big data are real data, such as automobile driving data and space observation data, generated from real world measurement and observation, social data derived from social media, e.g., Twitter and Instagram, and open data published by highly public groups, e.g., weather data and evacuation location data. These are generally called social

big data (SBD) [1]. Furthermore, SBD are inherently represented by multimedia (MM). By integrating and analyzing social big data, new knowledge can be obtained, which is expected to bring new value to society [2] [3].

SBD can include the same type of spatially different data, data obtained by different means for the same object, and temporally different data for the same object as well. Therefore, SBD applications cover not only use cases that include social data, but also use cases that include only engineering or scientific data generated in the real world.

Further as the horizon of applications whose main task is data analysis spreads, the following problems have emerged:

- *Application to science, e.g., lunar and planetary science*

Analytical applications in this field require strictness as science. That is, explanation of the protocol (procedure) of analysis and explanation of the reason for decisions are required [1]. In addition, as to the interpretation of the analytical model, it is necessary to explain the input data (for learning and test) and the data manipulation on the data, and the procedure (algorithm and program) for model construction. In order to interpret the individual results, it is necessary to explain the input data (actual data) and the reasons for the decisions.

- *Application to Social Infrastructure, e.g., Mobility as a Service (MaaS)*

Analytical applications in this field require consent of practitioners. That is, the analysis result must be consistent with the practitioners' own experiences, and especially in the case of applications such as ones related to human life, it is

necessary to fulfill the accountability to the concerned parties. Interpretation of both of the model and individual results is necessary as with science. In addition, especially if the data about the generic users are utilized in applications, interpretation of the model is also important in order to get rid of the general users' concerns.

In order for social big data to widely be used, it is necessary to explain the user the application system. Both microscopic description, that is, interpretation of the analytical model and explanation of individual decisions and macroscopic description, that is, description of the whole process including the data manipulation and the model construction are required.

First of all, the reason why a macro explanation is necessary is described below. In order for social big data applications to be accepted by users, it is necessary to ensure at least their reliability. Since information science is one area of science, we should guarantee reproducibility as science. In other words, it is necessary to ensure that third parties can prepare and analyze data according to given explanation and can get the same results.

In addition, in order for the service to be operable, it is necessary for the final user of the service to be convinced of how the service processes and uses the personal information. In addition, if the users can be convinced of the description of way of using the personal information, the progress of data portability can be advanced based on the EU's GDPR (General Data Protection Regulation) law on personal information protection [4] and Japan-based information bank to promote the use of personal information [5].

Next, a micro explanation is necessary for the following reasons. In order for analysts of social big data and field experts using the data to accept decisions made by the constructed model, it is assumed that they must understand the structure, actions and grounds of the model and are satisfied with them as well.

Up to now, the authors have been involved in the development of a wide range of social big data use cases ranging from tourism, disaster prevention to lunar and planetary science [6] [7]. In the course of these processes, from the users of the use cases, we have often received questions as to what kind of data are processed, what kind of models are created as the core of analysis, and furthermore, what are the grounds for the decisions. In other words, from the development experiences of multiple use cases, we have come to think that both the macro explanation proposed in this paper and the micro explanation emerging in AI are urgently needed.

To date, the authors created multiple seismic source classifiers of the lunar quakes (i.e., moonquakes) in the field of lunar and planetary science using the Balanced Random Forest [8], and the features, e.g., the distance between the moon and the earth, were calculated and studied for extracting features strongly related to the cause of

moonquakes as a micro explanation [6]. With regard to a macro explanation, the authors also showed that by observing many use cases, social big data applications should include different digital ecosystems such as data management (database operation) and data analysis (data mining, machine learning, artificial intelligence), we have noticed that it is necessary to have a method to generally describe the whole process of application consisting of such a hybrid digital ecosystem. Therefore, as a framework to describe processes in an abstraction level independent of a specific programming language, we have come to think of adopting a data model [9] developed in the field of database and proposed a framework for its description using the mathematical concept of family of sets [10]. As described in the subsequent section of the related works, the research on micro explanations is being actively carried out, whereas as far as research on the framework for a macroscopic description is not known except for our work.

The main contribution of the paper is to address the necessity of both macro and micro explanations for SBD applications and to propose an explanation framework integrating both of them. This will allow SBD applications to be more widely accepted and used. Although this paper describes our research-in-progress, we propose an integrated framework for explanation and introduce a part of its functions through case studies.

The contributions of this paper can be detailed as follows. First, as a science, a system that includes Artificial Intelligence (AI) or Data Mining (DM) needs reproducibility (How) [11] to ensure its reliability. For that purpose, it is important to show the procedures of the system explicitly and abstractly. We propose a dedicated data model for that purpose. For AI or DM, scientists need to show what model features are useful for making decisions and why the results obtained are valid (Why) [12]. This paper has scientific value in clarifying what features contribute more to the classification model and what can be shown as a basis for individual judgment through two use cases. The procedure can be modeled using a data model approach based on the mathematical concept *family* [10], using social data in the first case related to social science. The difference method [13] was used in the first case and the third case, related to natural science in order to model the hypotheses. In the explanation of the features of the analytical model in the second case, there is a skew in the data size for moonquake data, so we used Balanced Random Forest [8]. In the third case, for the basis of individual judgment we used CNN (deep learning) [14] and Grad-CAM (attention) [15] using Digital Elevation Model (DEM) provided by the Japan Aerospace Exploration Agency (JAXA).

This paper is of scientific value in that it demonstrates through use cases what can be explained to scientists as a basis for validating the results obtained by AI or DM. In other words, moonquake classification is important in lunar and

planetary science to understand the internal structure of the moon. This paper illustrated which features contribute most to the classification. The crater with a central hill is also a promising place for exploring the internal structure of the moon. This paper could illustrate what is the basis for judging the craters with central hills. These are also scientifically significant in that they have shown the possibility that AI and DM, which are IT technologies, are accepted by scientists as scientific methods.

The differences between this paper and the international conference paper [1] are as follows. The contribution and scientific value of this study were described in more detail. A description of the basic elements of the framework for explanation and the mechanism of its processing was added. A case of discovery of lunar craters with central hills was added as an example of a microscopic-explanation function (i.e., description of the grounds of judgment). The description of each use case was summarized according to the items such as scientific objectives, data, methods, and results.

In Section II, we will introduce our explanation framework. Through use case examples of macroscopic description and microscopic description, we will describe features of the proposed approach in Sections III, IV, and V, respectively.

II. OUR APPROACH

A. Explanation Framework

For a macro explanation of applications, the goal is to facilitate a data model for abstractly describing the entire processes from data acquisition to data analysis and to explain the processes based on the description. For the micro explanation, we aim to show the basis of the interpretation of the constructed model and the individual decisions made when applying it.

Macroscopic-explanation function (F1)
 = data management procedure + model generation procedure
Microscopic-explanation function
 = model feature explanation (F2) + judgement basis explanation (F3)

Figure 1. Explanation framework

The features of the proposed framework are summarized as follows.

Based on the SBD model introduced in Section III, the parties who are users of the framework (application developers) can describe the procedures (i.e., data management and model generation) of the application system in a more abstract manner than programming languages. The framework outputs the described procedures as they are as a macro explanation to the parties (e.g., tourism

operators in the tourism case). As a micro explanation, based on the results of the actual execution of the classification model, the framework outputs the features of the classification model (i.e., which features contribute to the classification) and the basis for judgment of each classification result to the scientists in the moonquake case and those in the moon crater case as the parties, respectively.

Figure 1 shows the framework. We specify which function corresponds to each use case. Case One in Section III illustrates the macroscopic-explanation function (F1) that explains the application procedures. Case Two in Section IV illustrates the microscopic-explanation function (F2) that explains which features contribute to the model classification, Case Three in Section V illustrates the microscopic-explanation function (F3) that explains the basis for individual judgment of the model classification.

We describe the framework in more detail as follows.

1) Construction of a theoretical foundation for integrated explanation

For that purpose, we build a theoretical framework of the technical foundation that integrates the following microscopic- and macroscopic-explanatory methods.

a) Macro explanation function: The application system is a hybrid ecosystem consisting of data management and data mining (including machine learning and Artificial Intelligence, or AI), and the function must be able to describe the application seamlessly. Moreover, it must be able to describe the application in a high level not depending on individual environments or programming languages. For instance, we aim to enable to describe “partition foreign visitors’ tweets into grids based on geo-tags.” Therefore, we first create a framework to unify the hybrid ecosystem based on the data model approach. In other words, we develop a method to provide macro explanations with the constituent elements (data structure and data manipulation) of the model based on the mathematical family of sets as a basic unit. The explanation mechanism provided by the proposed framework presents as a macro explanation a sequence of operations on databases to the user based on the model of SBD applications consisting of data management and data mining as in a use case depicted in Section III.

b) Microscopic-explanatory function: We develop an explanatory method independent of analytical model by extending explanatory functions based on attributes or constituent elements, which is an emergent approach in AI, discussed in the related work subsection. In other words, in model categories for structured data consisting of attributes, such as Support Vector Machine (SVM) and decision trees, we develop a method for systematically discovering subsets of attributes with strong influence on analysis results based on multiple weak classifiers. For instance, we aim to enable to illustrate a possibility that the features of the Earth and some of the features of Jupiter are effective for classification of the moonquakes when the moon is the origin of the

coordinate system. Especially this function is used to interpret the model itself. In model categories like Deep Neural Network (DNN) suitable for non-structured data such as images, we develop a method of explaining the analysis result based on the constituent elements or decomposition of the image with the use of annotation or attention. Especially this function is used to show the basis of individual decisions. For the micro explanation of the reasons for decisions, if the analysis target is image data, a part of the image which leads to the conclusion is indicated by concepts or words as its annotations based on a heat map. For instance, we aim to enable to illustrate that the contribution area for “central-peak crater” on the moon has heating area inside the crater, and the heating area is covering a central peak. If the object is structural data, that is, it consists of attributes, the micro explanation is presented in terms of the contribution ratios of the attributes as in a use case depicted in Section IV.

Please also note that data management and model construction in SBD applications are more complex than linear model construction frequent in traditional applications.

2) *Collection of use cases and verification of basic technology*

First, we collect several different kinds of use cases (tourism, mobility service, lunar exploration). We generate concrete explanations as targets for typical ones, using the integrated explanatory platform developed in items *a* and *b* and verify its feasibility

3) *Implementation of Explanation generation and presentation method*

Based on the theoretical framework of the integrated infrastructure, an automatic generation method of explanation and a presentation function of explanations are implemented. We evaluate their effectiveness by performing the experiments. We also incorporate InfoGraphics [16] as a method of presenting explanations to users since the users are not always analysis experts.

Basically, for micro explanation, we create explanations of individual decisions by solving partial problems that restrict information existing in original problems.

In this research, we aim to develop both the emerging microscopic-explanatory functions and macroscopic-explanatory functions and to build a framework for integrating two kinds of explanations.

B. Related Research

As a trend other than the authors' research, research corresponding to microscopic-explanatory functions has become active in AI, what is so called eXplainable AI (XAI) at present.

First, there is an attempt [17] to try to give a basic definition to the possibility of interpretation of a model in machine learning and research [18] on the evaluation method of interpretability.

Next, individual studies on XAI are roughly classified into (1) description based on features, (2) interpretable model, and (3) derivation of explanation model. Research is done to create a classification rule for explanation by creating a subset of features in SVM as a category of (1) [19]. In addition, in the image classification using Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM), there is research to generate explanations based on both image features and class features [20]. Further there is research introducing the explanation vector to make explicit the most important attributes [21]. In the category of (2), there is research using a AND/OR tree to discover the components of the model [22] and research to make models that can be interpreted by considering the generation process of features [23]. Research deriving description with reference of any classifier of the local approximation model falls into the category (3) [24].

In particular, the paper [25] is related to our framework. *Post-hoc global* explanation introduced in this paper corresponds to the explanation of the model features (micro explanation F2) in the proposed framework, *Post-hoc local* explanation introduced in the paper corresponds to the explanation of the basis for the judgement (micro explanation F3). However, the paper differs from our work in that the former takes no account of the macro explanation F1 (data management and model generation) in the proposed framework.

While developing along the approaches of (1) and (3) as a microscopic-explanatory technique, we aim to build a comprehensive explanation basis by conducting research on macroscopic-explanation technology.

In addition, although there is an application of infographics to a tourism use case [26], our research aims at basic research that can be widely used for visualization of explanation of general analysis.

III. CASE STUDY: MACRO EXPLANATION OF TOURISM APPLICATION

We will describe the case that explains how our data is used in analysis application. For that purpose, an integrated data model is introduced as a macroscopic description of an analytical application which is a hybrid ecosystem. Thus, the application is described using the integrated model just as a basis for macro explanation (see Figure 1). In other words, the data model introduced for model construction and reuse in the previous works [3] [13] is used for different purposes, i.e., the explanation functions for hypothesis generation.

A. Integrated Model

In the following subsections, we describe our data model approach to SBD, which consists of both of data structures and operations [9].

1) *Data model for SBD*

Our SBD model uses a mathematical concept of a *family*

[10], a collection of sets, as a basis for data structures. Family can be used as an apparatus for bridging the gaps between data management operations and data analysis operations.

Basically, our database is a *Family*. A Family is divided into *Indexed family* and *Non-Indexed family*. A Non-Indexed family is a collection of sets.

An Indexed family is defined as follows:

- a) $\{Set\}$ is a Non-Indexed family with Set as its element.
- b) $\{Set_i\}$ is an Indexed family with Set_i as its i -th element. Here, i : Index is called indexing set and i is an element of Index.
- c) Set is $\{<time space object>\}$.
- d) Set_i is $\{<time space object>\}_i$. Here, object is an identifier to arbitrary identifiable user-provided data, e.g., record, object, and multimedia data appearing in social big data. Time and space are universal keys across multiple sources of social big data.
- e) $\{Indexed\ family_i\}$ is also an Indexed family with Indexed family $_i$ as its i -th element. In other words, Indexed family can constitute a hierarchy of sets.

Please note that the following concepts are interchangeably used in this paper.

- Singleton family \Leftrightarrow set
- Singleton set \Leftrightarrow element

As described later in this section, we can often observe that SBD applications contain families as well as sets and they involve both data mining and data management. Please note that a family is also suitable for representing hierarchical structures inherent in time and locations as well as matrices and tensors associated with social big data.

If operations constructing a family out of a collection of sets and those deconstructing a family into a collection of sets are provided in addition to both family-dedicated and set-dedicated operations, SBD applications will be described in an integrated fashion by our proposed model.

2) SBD Operations.

SBD model constitutes an algebra with respect to Family as follows. SBD is consisted of Family data management operations and Family data mining operations. Further, Family data management operations are divided into Intra Family operations and Inter Family operations.

First, Intra Family Data Management Operations will be described as follows:

- a) *Intra Indexed Intersect* ($i:Index\ Db\ p(i)$) returns a singleton family (i.e., set) intersecting sets which satisfy the predicate $p(i)$. Database Db is an indexed Family, which will not be mentioned hereafter.
- b) *Intra Indexed Union* ($i:Index\ Db\ p(i)$) returns a singleton family union-ing sets which satisfy $p(i)$.

- c) *Intra Indexed Difference* ($i:Index\ Db\ p(i)$) returns a singleton family, that is, the first set (i.e., a set with smallest index) satisfying $p(i)$ minus all the rest of sets satisfying $p(i)$
- d) *Indexed Select* ($i:Index\ Db\ p(i)\ sp(i)$) returns an Indexed family with respect to i (preserved) where the element sets satisfy the predicate $p(i)$ and the elements of the selected sets satisfy the selection predicate $sp(i)$. As a special case of true as $p(i)$, this operation returns the whole indexed family. In a special case of a singleton family, Indexed Select is reduced to Select (a relational operation).
- e) *Indexed Project* ($i:Index\ Db\ p(i)\ a(i)$) returns an Indexed family where the element sets satisfy $p(i)$ and the elements of the sets are projected according to $a(i)$, attribute specification. This also extends also relational Project.
- f) *Intra Indexed cross product* ($i:Index\ Db\ p(i)$) returns a singleton family obtained by product-ing sets which satisfy $p(i)$. This is extension of Cartesian product, one of relational operators.
- g) *Intra Indexed Join* ($i:Index\ Db\ p(i)\ jp(i)$) returns a singleton family obtained by joining sets which satisfy $p(i)$ based on the join predicate $jp(i)$. This is extension of Join (a relational operator).
- h) *Select-Index* ($i:Index\ Db\ p(i)$) returns $i:Index$ of set $_i$ which satisfy $p(i)$. As a special case of true as $p(i)$, it returns all index.
- i) *Make-indexed family* (Index Non-Indexed Family) returns an indexed Family. This operator requires order-compatibility, that is, that i corresponds to i -th set of Non-Indexed Family.
- j) *Partition* ($i:Index\ Db\ p(i)$) returns an Indexed family. Partition makes an Indexed family out of a given set (i.e. singleton family either w/ or w/o index) by grouping elements with respect to p ($i:Index$). This is extension of "groupby" as a relational operator.
- k) *ApplyFunction* ($i:Index\ Db\ f(i)$) applies $f(i)$ to i -th set of DB, where $f(i)$ takes a set as a whole and gives another set including a singleton set (i.e., Aggregate function). This returns an indexed family. $f(i)$ can be defined by users.

Here the operations a) to g) are extensions of corresponding relational operators.

Second, Inter Family Data Management Operations will be described as follows:

All are assumed to be Index-Compatible.

- a) *Indexed Intersect* ($i:Index\ Db1\ Db2\ p(i)$) union-compatible

- b) *Indexed Union* ($i:Index\ Db1\ Db2\ p(i)$) union-compatible
- c) *Indexed Difference* ($i:Index\ Db1\ Db2\ p(i)$) union-compatible
- d) *Indexed Join* ($i:Index\ Db1\ Db2\ p1(i)\ p2(i)$)
- e) *Indexed cross product* ($i:Index\ Db1\ Db2\ p(i)$)

Finally, Family Data Mining Operations will be described as follows:

- a) *Cluster* (Family method similarity $\{par\}$) returns a Family as default, where Index is automatically produced. This is an unsupervised learner.
- b) *Make-classifier* ($i:Index\ set:Family\ learnMethod\ \{par\}$) returns a classifier (Classify) with its accuracy. This is a supervised learner.
- c) *Classify* (Index/class set) returns an indexed family with class as its index.
- d) *Make-frequent itemset* (Db supportMin) returns an Indexed Family as frequent itemsets, which satisfy supportMin.
- e) *Make-association-rule* (Db confidenceMin) creates association rules based on frequent itemsets Db, which satisfy confidenceMin. This is out of range of our algebra, too.

Please note that the predicates and functions used in the above operations can be defined by the users in addition to the system-defined ones such as Count.

B. Tourist Applications

First, we will summarize this case as follows.

- a) *(Social scientific and explanatory objectives)* In this case related to social science, it is important for the EBPM (Evidence-Based Policy Making) [27] parties (tourist operators) to identify where there is a gap between social needs (many foreigners want to use the Internet) and the infrastructure to meet them (free Wi-Fi access spots for foreigners are available). The procedure to realize this consists of data management and model generation (data mining and difference method). In other words, it is necessary to explain to the EBPM parties how to draw the conclusions (results of gap detection).
- b) *(Data used for use case)* Social media data Flickr [28] images and Twitter [29] articles were used. We collected 4.7 million Tweet articles (tweets) with geo tags by using the site provided API and selected 7,500 tweets posted by foreign visitors in Yokohama. We also collected 0.6 million Flickr images by using the site provided API and selected 2,100 images posted by

foreign visitors in Yokohama.

- c) *(Methods used for use case)* We used SQL to prepare social data and used a dedicated DM technique [30] to select only data posted by foreign visitors. We calculated the final result by using the difference method [3] [13] on separate results obtained from to the different data sources.
- d) *(Result)* As a result of social science, we could identify the areas with the gaps between social needs and available infrastructures. The model-based explanation of the whole processes for obtaining the result was found useful by talks with tourism operators.

Next, we will describe the case in more depth.

We describe a case study, finding candidate access spots for accessible Free Wi-Fi in Japan [31]. This case is classified as integrated analysis based on two kinds of social data.

This section describes our proposed method of detecting attractive tourist areas where users cannot connect to accessible Free Wi-Fi by using posts by foreign travelers on social media.

Our method uses differences in the characteristics of two types of social media:

Real-time: Immediate posts, e.g., Twitter

Batch-time: Data stored to devices for later posts, e.g., Flickr

Twitter users can only post tweets when they can connect devices to Wi-Fi or wired networks. Therefore, travelers can post tweets in areas with Free Wi-Fi for inbound tourism or when they have mobile communications. In other words, we can obtain only tweets with geo-tags posted by foreign travelers from such places. Therefore, areas where we can obtain huge numbers of tweets posted by foreign travelers are identified as places where they can connect to accessible Free Wi-Fi and /or that are attractive for them to sightsee.

Flickr users, on the other hand, take many photographs by using digital devices regardless of networks, but whether they can upload photographs on-site depends on the conditions of the network. As a result, almost all users can upload photographs after returning to their hotels or home countries. However, geo-tags annotated to photographs can indicate when they were taken. Therefore, although it is difficult to obtain detailed information (activities, destinations, or routes) on foreign travelers from Twitter, Flickr can be used to observe such information. In this study, we are based on our hypothesis of "A place that has a lot of Flickr posts, but few Twitter posts must have a critical lack of accessible Free Wi-Fi." We extracted areas that were tourist attractions for foreign travelers, but from which they could not connect to accessible Free Wi-Fi by using these characteristics of social media. What our method aims to find is places currently without accessible Free Wi-Fi.

Our method envisaged places that met the following two conditions as candidate access spots for accessible free Wi-Fi:

- *Spots where there was no accessible Free Wi-Fi*
- *Spots that many foreign visitors visited*

We use the number of photographs taken at locations to extract tourist spots. Many people might take photographs of subjects, such as landscapes based on their own interests. They might then upload those photographs to Flickr. These locations at which many photographs had been taken might also be interesting places for many other people to sightsee or visit. We have defined such places as tourist spots. We specifically examined the number of photographic locations to identify tourist spots to find locations where photographs had been taken by a lot of people. We mapped photographs that had a photographic location onto a two-dimensional grid based on the location at which a photograph had been taken to achieve this. Here, we created individual cells in a grid that was 30 square meters. Consequently, all cells in the grid that was obtained included photographs taken in a range. We then counted the number of users in each cell. We regarded cells with greater numbers of users than the threshold as tourist spots.

[Integrated Hypothesis] Based on different data generated from Twitter and Flickr, the following fragment as the macro explanation for hypothesis generation discovers attractive tourist spots for foreign visitors but without accessible free Wi-Fi currently (see Figure 2):

$DB_{t/visitor} \leftarrow$ Tweet DB of foreign visitors obtained by mining based on durations of their stays in Japan;

$DB_{f/visitor} \leftarrow$ Flickr photo DB of foreign visitors obtained by mining based on their habitations;

$T \leftarrow$ Partition (i : Index grid $DB_{t/visitor}$ $p(i)$); This partitions foreign visitors' tweets into grids based on geo-tags; This operation returns an indexed family.

$F \leftarrow$ Partition (j : Index grid $DB_{f/visitor}$ $p(j)$); This partitions foreign visitors' photos into grids based on geo-tags; This operation returns an indexed family.

$Index1 \leftarrow$ Select-Index (i : Index T $Density(i) \geq th1$); $Density$ counts the number of foreign visitors per grid. $th1$ is a threshold. This operation returns a singleton family.

$Index2 \leftarrow$ Select-Index (i : Index F $Density(i) \geq th2$); $Density$ also counts the number of foreign visitors per grid. $th2$ is a threshold. This operation returns a singleton family.

$Index3 \leftarrow$ Difference ($Index2$ $Index1$); This operation returns a singleton family.

Please note that Partition and Select-Index are family data management operations while Difference is a relational (set) data management operation.

We collected more than 4.7 million data items with geo-tags from July 1, 2014 to February 28, 2015 in Japan. We detected tweets tweeted by foreign visitors by using the method proposed by Saeki et al. [30]. The number of tweets that was tweeted by foreign visitors was more than 4.7

million. The number of tweets that was tweeted by foreign visitors in the Yokohama area was more than 7,500. We collected more than 0.6 million photos with geo-tags from July 1, 2014 to February 28, 2015 in Japan. We detected photos that had been posted by foreign visitors to Yokohama by using our proposed method. Foreign visitors posted 2,132 photos. For example, grids indexed by $Index3$ contain "Osanbashi Pier." Please note that the above description doesn't take unique users into consideration. The visual comparison of the same grids with unbalanced densities can help the decision makers to understand the proposal.

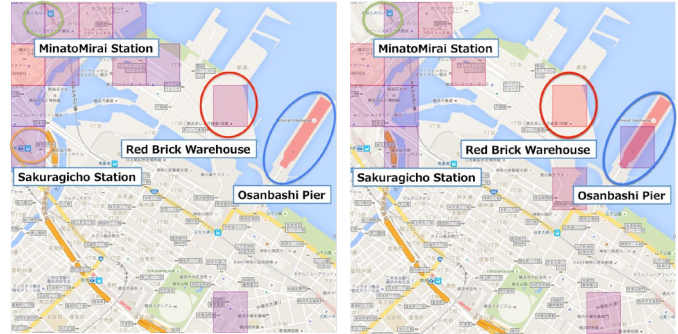


Figure 2. Differences of high-density areas of Tweets (left) and of Flickr photos (right).

IV. CASE STUDY: MICRO EXPLANATION FOR MOONQUAKE APPLICATION

First, we will summarize this case as follows.

- (Scientific and explanatory objectives) In this case related to lunar and planetary science, in order to know the internal structure of the moon, it is necessary to analyze the moonquake. As a preliminary study, the classification of moonquakes based on multiple epicenters is indispensable. However, it is not fully understood what features are more effective for moonquake classification. Therefore, it is necessary to determine the features that contribute most to the classification result as an explanation of the classification model.
- (Data used for use case) We used passive seismic data regarding to the moonquakes collected by the NASA Apollo program. The dataset [32] has 16 seismic sources and 2,480 events as depicted in Table II. There is a skew with respect to the size of each source.
- (Methods used for use case) We used Balanced Random Forest [8] to explain which features most contribute to one-to-one classification of moonquakes with respect to seismic sources.
- (Result) Results of the classification performance using orbit parameters of objects in our Solar System (Earth, Sun, Jupiter, and Venus) suggest that the

Earth orbit parameter is the most effective feature among them. The Jupiter orbit parameter is effective for classification of some seismic sources. The effect was validated by discussions in our research team consisting of IT specialist and natural scientists.

Next, we will describe the case of determining features important for interpreting the constructed model by reducing features with small contribution ratios. We apply Balanced Random Forest [8], which extends Random Forest [33], a popular supervised learning method in machine learning, to lunar and planetary science to verify the key features in analysis. Our verification method tries to confirm whether the known seismic source labels can be reproduced by Balanced Random Forest using the features described below based on the features constructed from the moonquakes with the seismic source label of the known moonquake as the correct label.

A. Features for Analysis

TABLE I shows the parameters in the coordinate systems used in this section. We use as seismic source of moonquakes the position on the planets of the moon, the sun, the earth, and Jupiter (X, y, z), velocity (v_x, v_y, v_z), and distance (lt). Based on the time of moonquake occurrence, we calculate and use features using SPICE [34]. SPICE assists scientists in planning and interpreting scientific observations from space-borne instruments, and to assist engineers involved in modeling, planning and executing activities needed to conduct planetary exploration missions.

Here, sun perturbation is the solar perturbation. The IAU MOON coordinate system is a fixed coordinate system centered on the moon. The z axis is the north pole direction of the moon, the x axis is the meridian direction of the moon, the y axis is the right direction with respect to the plane xz. The IAU EARTH coordinate system is a fixed coordinate system centered on the earth. Here, the z axis is the direction of the conventional international origin, the x axis is the direction of the prime meridian, and the y axis is the right direction with respect to the xz plane.

We also calculate the period of the perigee at the distance of *earth_from_moon*, the period based on the period of the perigee, the periods of the x coordinate and the y coordinate of the solar perturbation. *sin* and *cos* values are calculated from these periodic features and the phase angle based on them. In addition, at the positions *moon_from_earth* and *sun_from_earth*, we calculate the *cos* similarity as the features of the sidereal moon. Most importantly, as all possible combinations of these features, a total of 55 features are used in our experiments described here.

B. Balanced Random Forest

Random Forest is an ensemble learning that combines a large number of decision trees and is widely used in fields such as data mining and has a characteristic that the

contribution ratio of features can be calculated. However, Random Forest has a problem such that when there is a large difference in the size of data to be learned depending on class labels, the classifier is learned biased towards classes with a large size of data. Generally, we address the problem of imbalanced data by weighting classes with a small number of data. However, if there is any large skew between the numbers of data, the weight of data belonging to classes with a small number will become large, which is considered to cause over fitting to classes with a small number of data. Since the deep moonquakes have a large difference in the number of events for each seismic source, it is necessary to apply a method considering imbalanced data.

As analysis considering imbalanced data, we apply Balanced Random Forest [8], which makes the number of samples even for each class when constructing each decision tree. Balanced Random Forest divides each decision tree based on the Gini coefficient. Gini coefficient is an index representing impurity degree, which takes a value between 0 and 1. The closer it is to 0, the higher the purity is, that is, the less variance the data have. The contribution ratio of the feature is calculated for each feature by calculating the reduction ratio by the Gini coefficient at the branch of the tree. The final contribution ratio is the average value of contribution ratios of each decision tree.

C. Experiment Setting

Here, we describe experiments for evaluating features effective for seismic source classification, together with the results and considerations. Based on the classification performance and the contribution ratio of the features by Balanced Random Forest, we analyze the relationship between the seismic sources in the features used in this paper.

The outline of feature analysis is as follows: Features are calculated based on the time of occurrence of moonquake. Balanced Random Forest is applied to each pair of all seismic sources. Classification performance and the contribution ratio of the features by Balanced Random Forest are calculated and analyzed.

In this paper, as one-vs-one method, by constructing the classifier for every pair of two seismic sources in the dataset, we perform analysis paying attention to characteristics of each seismic source and the relationship between seismic sources. 100 Random Forests are constructed for each classifier. The number of samples used to construct each decision tree are taken 50 by bootstrap method. Bootstrap is a test or metric that relies on random sampling with replacement. Also, scikit-learn [35] was used to construct each decision tree in Random Forest. scikit-learn is a machine learning library for the Python programming language. In this paper, we perform the following analysis as feature selection.

- We create a classifier that learns all of the extracted 55 features.

TABLE I. PARAMETERS IN THE COORDINATE SYSTEMS COMPUTED USING SPICE.

Target	Observer	Coordinate system	Parameter
EARTH BARYCENTER	MOON	IAU MOON	earth_from_moon
SOLAR SYSTEM BARYCENTER	MOON	IAU MOON	sun_from_moon
JUPITER BARYCENTER	MOON	IAU MOON	jupiter_from_moon
SOLAR SYSTEM BARYCENTER	EARTH BARYCENTER	IAU EARTH	sun_from_earth
JUPITER BARYCENTER	EARTH BARYCENTER	IAU EARTH	jupiter_from_earth
SUN	SOLAR SYSTEM BARYCENTER	IAU EARTH	sun_perturbation

TABLE II. NUMBER OF DATA FOR EACH SEISMIC SOURCE.

Seismic source	A1	A5	A6	A7	A8	A9	A10	A14	A18	A20	A23	A25	A35	A44	A204	A218
Number of data	441	76	178	85	327	145	230	165	214	153	79	72	70	86	85	74

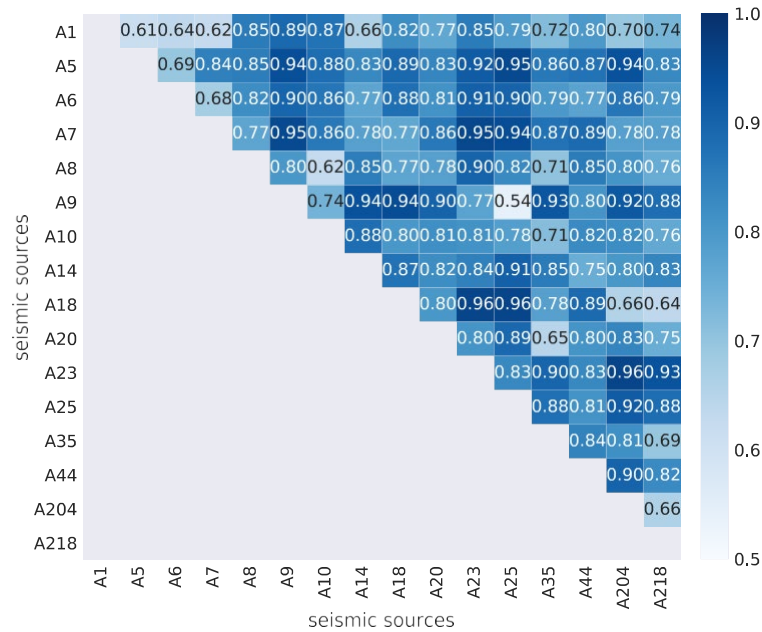


Figure 3. Averages of F-measures for pairs of seismic sources.

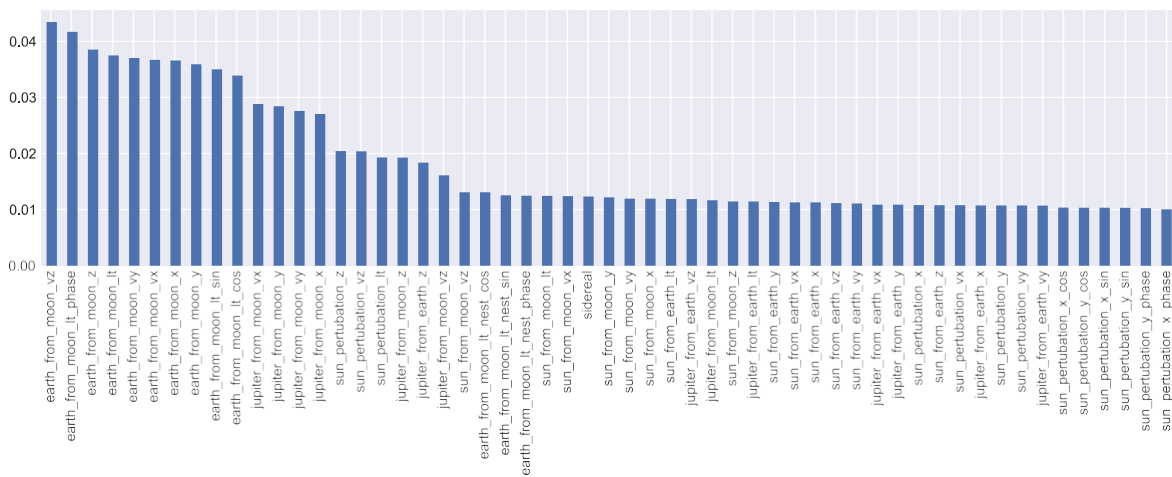


Figure 4. Averages of contribution ratios for each feature.

- Using the Variance Inflation Factor (VIF), we construct a classifier after reducing features. VIF quantifies the severity of multicollinearity, that is, a phenomenon in which one predictor variable in a multiple regression model can be linearly predicted from the others accurately.

Here, VIF is one of the indicators used to evaluate *multicollinearity*. In this paper, in order to make VIF of each feature 6 or less, experiments were conducted on a subset with reduced features. Based on the experimental results using all features, we calculate VIF and delete features with 6 or more VIF. To calculate VIF, statsmodel [36] was used.

TABLE II shows the dataset in this paper. We select events of 16 seismic sources whose observed number of moonquake events is 70 or more.

In this paper, the precision ratio, recall ratio, and F-value are used as indexes for evaluating the performance of classification of seismic sources.

The precision ratio is an index for measuring the accuracy of the classification, and the recall ratio is an index for measuring the coverage of the classification. F-value is the harmonic mean of recall and precision ratios and is an index in consideration of the balance of precision and recall. The score of the classifier in this paper is the average value of the F-values of the two classes targeted by the classifier.

D. Experiment Results

1) Experimental results using all features

a) Classification performance

Figure 3 is the average of the F-measures of classifiers for each seismic source. F-measure is the harmonic mean of precision and recall in statistical analysis. The vertical axis and the horizontal axis show seismic sources, each value is a score of the average of F-measure of classifier. In Figure 3, the highest classification performance is 0.96 and it is observed in multiple pairs of seismic sources. Also, the lowest classification performance is 0.54 as of classifier between A9 and A25. Figure 3 shows that some classification is difficult depending on combinations of seismic sources. Also, the number of classifiers with 0.9 or higher as classification performance is 20, about 17% of the total number of the classifiers. The number of classifiers with 0.8 or more and less than 0.9 is 60, 50% of the total. The number of classifiers with performance below 0.6 is only one. Most of the classifiers show high classification performance and show that the positional relationships of the planets are effective for the seismic source classification of the deep moonquakes.

b) Contribution ratio of features

Figure 4 shows the average value of contribution ratios for each feature. All features with the higher contribution ratios are those of the earth when they are calculated as the moon as the origin of the coordinate system. In addition, it shows that the contribution ratios of Jupiter's features are high when the moon is the origin. By comparing features when the moon is

the origin and when the earth is the origin, the features with the moon as the origin has a higher contribution ratio than the features with the earth as the origin. These observations suggest that the tidal forces are among the causes of moonquakes. Figure 4 indicates that relationships between the moon and the Earth affect the classification most strongly. However, there is a possibility that correlation between features, then it is necessary to further analyze each feature from view point of mutual independence. Therefore, in the following subsection, considering the correlations between features, we will describe the experimental results after feature reduction using VIF.

2) Experimental results of feature reduction using VIF.

a) Classification performance

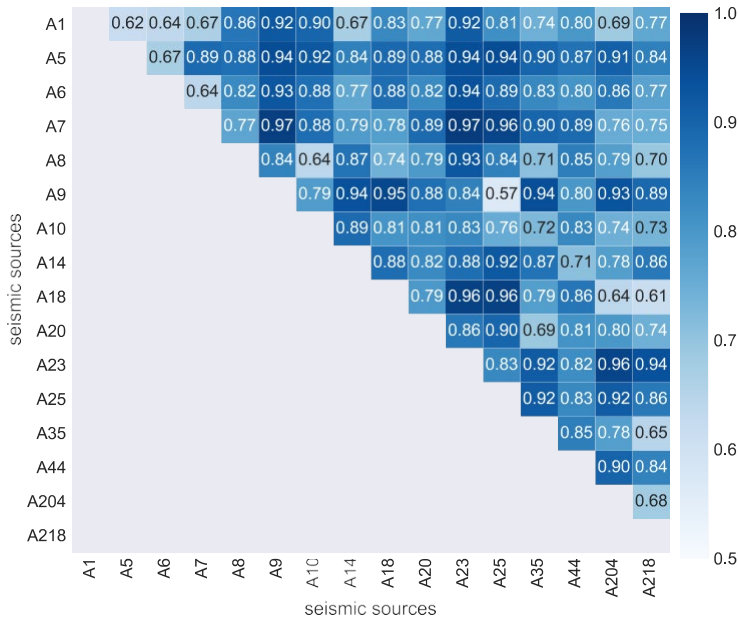
Figure 5 shows the average of the F-measures of the classifier when the features are reduced. Similarly, as in Figure 3, the vertical axis and the horizontal axis are seismic sources, respectively, and each value is the score of the F-measure of the classifier in Figure 5. In addition, the number of classifiers whose classification performance is 0.9 or higher is 26, about 22% of the total. 54 classifiers with 0.8 or higher but less than 0.9 are 45% of the total. There is one classifier whose classification performance is less than 0.6. Compared with Figure 3, these show that the classification performance does not change significantly.

b) Contribution ratio of features

Figure 6 shows the average value of the contribution ratios of each seismic source after feature reduction. After reducing features, earth features when the origin is the moon are reduced to 4 features of the top 10 features which existed before feature reduction. The four features between top 11 and 14 positions of the features of Jupiter when the origin is the moon, as shown in Figure 4, are reduced to one feature. Other parameters of Jupiter are thought to have been affected by other features. The subset of the features after feature reduction is considered to have small influence of multicollinearity. Therefore, there is a possibility that the features of the Earth and some of the features of Jupiter are effective for classification when the moon is the origin. These results are microscopic explanations made directly from the model constructed by Balanced Random Forest.

E. Discussion of methods and features

By using Balanced Random Forest, contribution ratios of features can be easily calculated in addition to classification performance, so it is useful for feature analysis like the scientific research described in this section. However, in this method, there is room for consideration of parameters of classification techniques depending on the seismic sources as the classification targets. Moreover, in order to obtain higher classification performance, it is necessary to consider many classification methods. Furthermore, it is necessary to apply a method considering waveform information simultaneously collected by the NASA Apollo project. In addition, since the



findings obtained in this paper are only correlations, it is difficult to directly estimate the causal mechanism of the deep moonquakes. However, the results of this paper are shown to be useful for further analysis and knowledge creation by experts. If the knowledge of experts such as the physical mechanism about moonquakes is available, the elucidation of the causal relationships between the seismic sources and the planetary bodies and ultimately that of the causal mechanism of the moonquakes (possibly related to tidal forces) can be expected. In general, expertise in any domain is expected to increase our understanding of the causal relationships suggested by our correlation analysis.

Figure 5. Averages of F-measures for pairs of seismic sources after feature reduction.

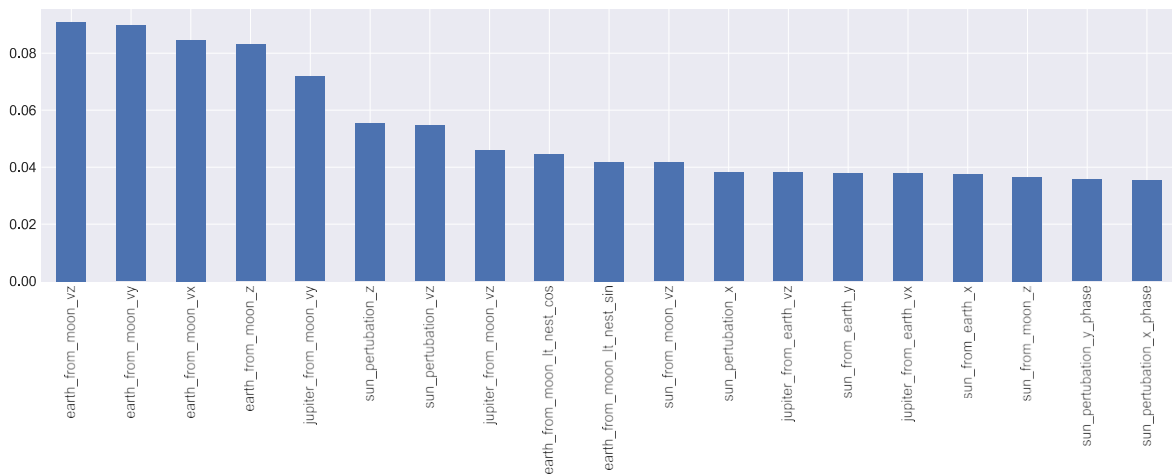


Figure 6. Averages of contribution ratios for each feature after feature reduction.

V. CASE STUDY: MICRO EXPLANATION FOR CENTRAL PEAK CRATER APPLICATION

First, we will summarize this case as follows.

- a) (Scientific and explanatory objectives) In this case, also related to lunar and planetary science, in order to understand the internal structure and movement of the moon, it is conceivable to use the materials inside the moon as a clue. The central hill in the crater is attracting attention as a place where the materials inside the moon are exposed on the moon

surface. However, not all craters with central hills on the moon have been identified. Therefore, it is scientifically necessary to make the catalog. Therefore, it is also necessary to explain to the relevant scientists the grounds for judging the craters included in the found candidates as craters with central hills.

- b) (Data used for use case) We used about 7,200 images provided by both NASA and JAXA. Each image has been resized to 512 (height) * 512 (width) * 1

(normalized elevation). We divided the whole images into equal numbers of images with three labels, that is, craters with central hills (central-hill craters), craters without central hills (normal craters), and non-craters.

- c) (Methods used for use case) We used RSPD [37] to detect craters only for preparation of training data of CNN [14][38]. Next, we applied learnt a CNN to find central-hill craters including unknown and known ones and used Grad-CAM [15] to know the evidence for judging central-hill craters.
- d) (Result) We could classify three classes with 96.9 % accuracy, which was verified by the scientific members of our research team. We also could show the scientific members of our research team individual evidences for judging central-hill craters consisting of crater rims and central hills.

Next, we will describe this case in more depth for the explanation functions although we have already introduced it to explain the way of model construction in our previous work [3].

A. Discovery of central peak craters

Scientific data are a kind of real-world data. By taking an example of research conducted by our team including JAXA researchers using scientific data which are also open data, we explain integrated analysis [39]. We use hypothesis generation based on differences between original data and their rotations.

A detailed map of the surface of the moon was provided by JAXA launched lunar orbiter KAGUYA (SELENE) [40]. Of course, KAGUYA's purpose goes beyond making a map of the moon. The goal is to collect data that will help elucidate the origin and evolution of the moon. In order to further pursue such purposes, it is important to examine the internal structure of the moon.

One of the methods to examine the internal structure of the moon is to analyze the data of moonquakes (i.e., corresponding to earthquakes) that occur in the moon. Research is also being conducted to classify the hypocenters of moonquakes based on the data of moonquakes provided by the NASA Apollo program. Among these studies are our research which showed that it is possible to classify moonquakes by features such as the distance between the moon and the planet alone without using seismic waves themselves as described in Section IV.

Another method is to launch a spacecraft to directly explore the internal structure of the moon. However, it is not sufficient to land the spacecraft anywhere on the moon. That is naturally because there are limited resources such as budgets that can be used for lunar exploration. In other words, it is necessary to determine the effective point as the target of the spacecraft based on the evidence. This way is an example of

EBPM, making an effective plan based on evidence under limited resources.

On the other hand, whether large or small, a lot of craters exist in the moon. Among them, special crater with a structure called "central peak" (hereinafter referred to as "central peak craters") is present (see Figure 7). The central peak is exposed on the moon and lunar crustal substances are also exposed therein. Therefore, it is likely that central peak craters scientifically have important features. In other words, the exploration of the surface of the central peak makes it possible to analyze the surrounding internal crustal materials in a relatively easy way. By this, it is expected that not only the origin of the crater and central peak can be estimated, but also the surface environment of the past lunar surface and the process of crustal deformation of the moon can be estimated.

But with respect to the central peak crater as the exploration target, conventionally the confirmation of existence of the central peaks has been visually done by the experts. So, the number of craters known as central peak craters is rather small. This problem can be solved by automatic discovery and cataloguing of central peak craters to significantly increase the number of central peak craters as candidate exploration points.

Thus, in this case with creating the catalog of central peak craters as our final goal, a specific technique for automatic discovery of central peak craters has been proposed. This case uses DEM (Digital Elevation Model) of the lunar surface as results observed by the lunar orbit satellite "KAGUYA" of JAXA [40]. Paying attention to the image characteristics of DEM, we apply CNN (Convolutional Neural Network [14] [38]) as a popular technique for deep learning, which is recently in the limelight as AI, to construct the discrimination model. We evaluate discriminability of the central peak crater by the model by experiments.

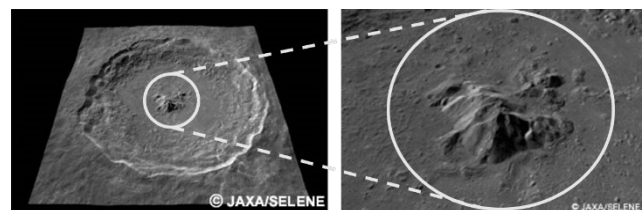


Figure 7. Example of central peak crater.

B. Integration Hypothesis

The central peak crater is identified by the following two-step procedure.

- 1) Crater extraction on the moon by RPSD method

Craters are extracted by using the popular and secure method called RPSD (Rotational Pixel Swapping for DTM) for digital terrain models. Here, DTM (Digital Terrain Model) is a digital terrain model similar to the digital elevation model (DEM). The RPSD method focusses on the rotational symmetry when rotating the image of the DEM at a certain point (i.e., central point). That is, RPSD uses the fact that the negative gradient property from the rim of the crater to the center in the lunar DEM does not change with rotation of craters. In other words, we make the difference between the original candidate crater and the rotated one (corresponding to the difference in the observation mode for the same object) and confirm that the feature about rotational symmetry does not change in order to discriminate craters. In a word, this method corresponds to our generalized difference method in which hypotheses (craters) are found by focusing on differences obtained by different means for the same object (candidate craters).

2) CNN-based automatic discrimination of central peak crater from extracted craters

In general, in the discrimination phase for each layer of deep learning, each output node multiplies the input values by weights, takes their sum, and adds their bias to the sum, and then outputs the result in the forward direction.

In the learning phase of deep learning, as a problem of minimizing the error between the output of discrimination and the correct answer, the values of weights and biases are updated by differentiating the error function with respect to the weight and bias of each layer.

C. Integrated Analysis

First, using RPSD, we extract the DEM data of each candidate crater and provide them with a label (non-crater, non-central peak crater, and central peak crater) to create training data. We learn CNN model thus using the training data and discriminate the central peak craters by using the CNN model. From recall ratios obtained by experiments focusing on how much correct answers are contained in the results, the possibility that CNN is an effective technique in the central peak crater determination is confirmed.

In order to confirm reasons for the classification results, we visualize the contribution areas in input images which affect the model (i.e., individual evidence).

We use Grad-CAM [15], a method for visualizing contribution areas for each label in an input image. We use it because it has an affinity for CNN.

The left part (see Figure 8) is an input image, the central part is the contribution area for “central-peak crater” label, and the right part is the contribution area for “normal crater” label.

The contribution area for “central-peak crater” has heating area inside the crater, and heating area is covering a central peak. On the other hand, the central peak area

does not have heating area at the contribution area for “normal crater.” Therefore, we think that the central peak area contributes classification for “central-peak crater” label. Thus, the contribution areas as the micro explanation can help the scientists to understand the corresponding classification results.

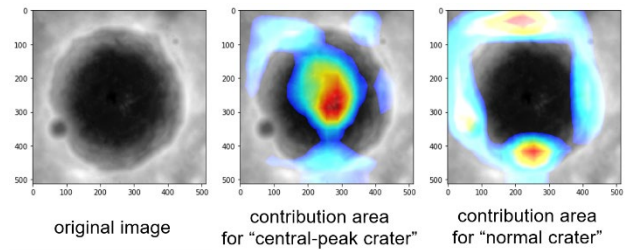


Figure 8. Explanation of individual evidences.

VI. CONCLUSION

In this paper, we proposed a general framework of explanation necessary to widely promote implementation of analytical applications using SBD. The procedure of a tourism application based on integrated data model was described as an example of a macro explanatory function. In addition, we used Balanced Random Forest as a micro explanatory function to extract features effective for the seismic source classification of the deep moonquakes from the temporal and spatial features of the planets. We described another example of a microscopic-explanatory function to explain individual evidences for discrimination of central peak craters.

The results of social science research (i.e., an example of macro explanation of procedures) were explained to external travel experts to confirm their effectiveness. Regarding to the explanation of scientific results (i.e., examples of micro explanations of model features and judgements), we have positive feedbacks from the relevant scientists in our research team based on the scientific effectiveness.

We reiterate the whole process of the SBD application with explanations as the summarization of the contribution of this paper.

1. The user describes the procedure for data management and model generation by utilizing the data model (i.e., SBD model) and the hypothesis generation methods (e.g., generalized difference method).
2. The macroscopic-explanation function uses the constructed description for the explanation.
3. The microscopic-explanation function finds the effective model features and individual judgement basis by executing the constructed model using the explanation-oriented techniques (e.g., Balanced Random Forest and Grad-CAM).

We will continue to verify the versatility of the explanatory framework by applying it to a wider variety of use cases in the future. We will also continue to interview the parties concerned and listen to the opinions of experts at international conferences on the effectiveness of the framework for explanation. In fact, we have already presented the candidate list of central hill craters with the micro explanations to the scientists in the related field. They have definitely found unidentified central hill craters among the candidates. As a result, a new project has recently been initiated to re-estimate the quantitative relationships holding between the radius of the central hill crater and the height of the central hill based on our findings.

ACKNOWLEDGMENT

This work was supported by JSPS KAKENHI Grant Number 16K00157, 16K16158, 20K12081 and Tokyo Metropolitan University Grant-in-Aid for Research on Priority Areas Research on social big data.

REFERENCES

- [1] H. Ishikawa, Y. Yamamoto, M. Hirota, and M. Endo, "Towards Construction of an Explanation Framework for Whole Processes of Data Analysis Applications: Concepts and Use Cases," Proc. The Eleventh International Conference on Advances in Multimedia MMEDIA 2019 (Special tracks:SBDMM: Social Big Data in Multimedia), Mar. 2019.
- [2] H. Ishikawa, Social Big Data Mining, CRC Press, 2015.
- [3] H. Ishikawa, D. Kato, M. Endo, and M. Hirota, "Applications of Generalized Difference Method for Hypothesis Generation to Social Big Data in Concept and Real Spaces," Proc. the 11th International Conference on Management of Digital EcoSystems (MEDES 2019), pp. 44–55, November 2019. <https://doi.org/10.1145/3297662.3365822> [retrieved: May, 2020].
- [4] EU GDPR, <https://eugdpr.org/> [retrieved: May, 2020].
- [5] J. Hemmi, "Japan's 'information banks' to let users cash in on personal data," NIKKEI Asian Review, May 09, 2019.
- [6] K. Kato, R. Yamada, Y. Yamamoto, M. Hirota, S. Yokoyama, and H. Ishikawa, "Analysis of Spatial and Temporal Features to Classify the Deep Moonquake Sources Using Balanced Random Forest," Proc. The Ninth International Conferences on Advances in Multimedia (MMEDIA 2017), April 2017.
- [7] T. Tsuchida, D. Kato, M. Endo, M. Hirota, T. Araki, and H. Ishikawa, "Analyzing Relationship of Words Using Biased LexRank from Geotagged Tweets," Proc. ACM MEDES International Conference, pp. 42-49, 2017.
- [8] C. Chen, A. Liaw, and L. Breiman, "Using random forest to learn imbalanced data," University of California, Berkeley, pp. 1-12, 2004.
- [9] H. Ishikawa, "Object-oriented database systems," (C. T. Leondes ed.) Database and Data Communication Network Systems Techniques and applications, vol. 1, pp. 77-122, Academic Press 2002.
- [10] D. Smith, R. St. Andre, and M. Eggen, A Transition to Advanced Mathematics, Brooks/Cole Pub Co., 2014.
- [11] R. D. Peng, "Reproducible Research in Computational Science," Science, vol. 334, issue 6060, pp. 1226-1227, 2011.
- [12] K. McCain, "Explanation and the Nature of Scientific Knowledge," Sci & Educ, vol. 24, pp. 827–854 (2015). <https://doi.org/10.1007/s11191-015-9775-5> [retrieved: May, 2020].
- [13] H. Ishikawa, D. Kato, M. Endo, and M. Hirota, "Generalized Difference Method for Generating Integrated Hypotheses in Social Big Data," Proc. ACM MEDES International Conference, pp. 13-22, 2018.
- [14] X.-W. Chen and X. Lin, "Big Data Deep Learning: Challenges and Perspectives," IEEE Access, vol. 2, pp. 514-525, 2014.
- [15] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization," Proc. The IEEE International Conference on Computer Vision (ICCV), 2017, pp. 618-626 <https://arxiv.org/abs/1610.02391> [retrieved: May, 2020].
- [16] W. V. Siricharoen, "Infographics: The New Communication Tools in Digital Age," Proc. International Conference on E-Technologies and Business on the Web, pp. 169-174, 2013.
- [17] Z. C. Lipton, "The Mythos of Model Interpretability," Communications of the ACM, vol. 61, no. 10, pp. 36-43, October 2018.
- [18] F. D. Velez and B. Kim, "A roadmap for a rigorous science of interpretability," pp. 1-13, 2017 (arXiv: 1702.08608, 2017).
- [19] D. Martens, B. Baesens, T. V. Gestel, and J. Vanthienen, "Comprehensible credit scoring models using rule extraction from support vector machines," Rule extraction from support vector machines, pp. 33-63, 2008.
- [20] L. A. Hendricks, Z. Akata, M. Rohrbach, J. Donahue, B.

- Schiele, and T. Darrell, "Generating visual explanations," Proc. European Conference on Computer Vision, pp. 3-19, Springer, 2016.
- [21] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Mueller "How to explain individual classification decisions," The Journal of Machine Learning Research, vol. 11, pp. 1803-1831, August 2010.
- [22] Z. Si and S. C. Zhu., "Learning and-or templates for object recognition and detection," IEEE Trans. Pattern Anal. Mach. Intell., vol. 35, no. 9, pp. 2189-2205, 2013.
- [23] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," Science, vol. 350, issue 6266, pp. 1332-1338, 2015.
- [24] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why Should I Trust You?: Explaining the Predictions of Any Classifier," Proc. CHI 2016 Workshop on Human Centered Machine Learning, pp. 1135-1144, 2016 (arXiv: 1602.04938v1 [cs.LG] 16 Feb 2016).
- [25] M. Du, N. Liu, and X. Hu, "Techniques for Interpretable Machine Learning," Communications of the ACM, vol. 63, no. 1, pp. 68-77, 2020.
- [26] K. W. Su, C. L. Liu, and Y. W. Wang, "A principle of designing infographic for visualization representation of tourism social big data." Journal Ambient Intell. Human Comput, pp. 1-21, 2018, doi:10.1007/s12652-018-1104-9.
- [27] OECD, "EBPM, Evidence in Education: Linking Research and Policy." <http://www.oecd.org/education/ceri/38797034.pdf> [retrieved: May, 2020].
- [28] Flickr <https://www.flickr.com/> [retrieved: May, 2020].
- [29] Twitter <https://twitter.com/> [retrieved: May, 2020].
- [30] M. Hirota, K. Saeki, Y. Ehara, and H. Ishikawa, "Live or Stay?: Classifying Twitter Users into Residents and Visitors," Proc. International Conference on Knowledge Engineering and Semantic Web (KESW 2016), pp. 1-2, 2016.
- [31] K. Mitomi, M. Endo, M. Hirota, S. Yokoyama, Y. Shoji, and H. Ishikawa, "How to Find Accessible Free Wi-Fi at Tourist Spots in Japan," Volume 10046 of Lecture Notes in Computer Science, pp. 389-403, 2016.
- [32] DARTS Available: <http://darts.jaxa.jp> [retrieved: May, 2020]
- [33] L. Breiman, "Random forests," Machine learning, vol. 45, no. 1, pp. 5-32, 2001.
- [34] NAIF, <https://naif.jpl.nasa.gov/naif/> [retrieved: March, 2019].
- [35] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, A. Müller, J. Nothman, G. Louppe, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay, "Scikit-learn: Machine learning in Python," Journal of Machine Learning Research, vol. 12, pp. 2825-2830, 2011.
- [36] S. Seabold and J. Perktold. "Statsmodels: Econometric and statistical modeling with python," Proc. 9th Python in Science Conference, pp. 57-61, 2010.
- [37] S. Yamamoto, T. Matsunaga, R. Nakamura, Y. Sekine, N. Hirata, and Y. Yamaguchi, "An Automated Method for Crater Counting from Digital Terrain Model Using Rotational Pixel Swapping Method," Proc. 49th Lunar and Planetary Science Conference, 2 pages, 2018.
- [38] Y. LeCun, Y. Bengio, and Geoffrey Hinton, "Deep learning," vol. 521, pp. 436-444, 2015. doi:10.1038/nature14539
- [39] S. Hara H. Inoue, M. Yamamoto, Y. Yamamoto, M. Otake, H. Otake, T. Araki, M. Hirota, and H. Ishikawa, "Automatic Extraction of Lunar Central Peak Craters by Deep Learning," 16th AOGS Annual Meeting, 2019.
- [40] JAXA, KAGUY(SELENE) Data Archive <https://darts.isas.jaxa.jp/planet/pdap/selene/> [retrieved: May, 2020].