DATA MINING SERVICE SEMANTICS IN CLOUD COMPUTING

#1 J. KUMARI, #2 N. MADHUSUDHANA RAO
#1 ASSISTANT PROFESSOR#2 MCA SCHOLAR
DEPARTMENT OF MASTER OF COMPUTER APPLICATIONS
QIS COLLEGE OF ENGINEERING & TECHNOLOGY, ONGOLE
VENGAMUKKAPALEM(V),ONGOLE,PRAKASAM DIST.,ANDHRA PRADESH

ABSTRACT

The recent integration of new Data Mining and Machine Learning services by Cloud Computing providers is equipping customers with very sophisticated data analysis tools that encompass all the benefits of this environment. Cloud computing service providers for data mining frequently offer descriptions and specifications in various forms that are often incompatible with those of other providers. From a functional perspective, the ability to delineate comprehensive Data Mining services is essential for preserving usability and, notably, the portability of these services, irrespective of software or hardware support, as well as variations among cloud platforms. The primary aim of this paper is to formulate a Data Mining service definition that enables the composition of a comprehensive service through a singular, straightforward definition, thereby facilitating the portability and deployment of a data mining workflow across various providers or within a marketplace of readily consumable services. This article introduces a semantic framework for defining and describing comprehensive Data Mining services, taking into account both the provider's administration of the service (pricing, authentication, Service Level Agreement, etc.) and the delineation of the Data Mining workflow as a service. It constitutes a significant advancement towards the standardization and industrialization of Data Mining services. To evaluate the scheme's validity, a catalog of services from Data Mining providers has been delineated, and a comprehensive example of a service for a Random Forest method has been established. A practical scenario has been established, building a deployment platform for Data Mining services to provide functional support to the scheme, demonstrating the practical advantages of the concept for the end user.

03779254 Page 192 of 203

I. INTRODUCTION

Cloud Computing has been seamlessly integrated into our daily life. The accessibility of the Internet and the rapid Proliferation linked of gadgets has significantly enhanced its popularity. Embracing the phenomenon of Cloud Computing signifies a major transformation in the exploration, consumption, deployment of Information Technology services. CCis a service delivery mechanism for enterprises, entities, and customers, adhering to the utility paradigm, similar to electricity or gas services. Cloud computing is a paradigm of service delivery in which computational resources and processing power are acquired via the Internet of Services. The volume of data generated by enterprises and organizations is escalating at an exceptionally rapid pace. Forbes forecasts that in 2020, growth will persist, with data generation anticipated to escalate by as much as 4,300%, driven by the substantial volume of data produced by service users. By 2020, it is projected that over 25 billion devices will be connected to the Internet, as per Gartner, generating in excess of 44 billion gigabytes of data each year. In this context, cloud computing providers are utilizing their extensive computational capacity to offer cloud clients

novel services for Data Mining. Cloud companies and services, including Amazon SageMaker and Microsoft Azure Machine Learning Studio, provide a collection of algorithms as services within computing platforms. In this context, other cloud computing platforms, such Algorithmic and Google Cloud ML, offer advanced Machine Learning capabilities, including object detection in photos, sentiment analysis, text mining, forecasting, among others. Each CC service provider possesses a distinct definition of these services. which is typically incongruent with those of other providers, affecting not only DM-related issues but also the administration of the CC service. For example, one supplier may offer a service utilizing a Random Forest algorithm, while another provider may present a similarly named method with comparable features or parameters, despite the two potentially being identical. This complicates the definition of services or service models independently of the provider and hinders the comparison of services via a cloud computing service broker. A standardization of service definitions would enhance competitiveness, enabling third parties to engage with these services transparently, bypassing the specific details of particular

03779254 Page 193 of 203

providers. Numerous suggestions exist for defining services, encompassing significant range of both syntactic and semantic languages to provide accurate specification and modeling of services. Proposals derived from Linked Data can address the issue of service definition from a more holistic perspective. LD employs models and frameworks from the Semantic Web, a technology designed to present data on the web in a more reusable and interoperable manner with other applications. The LD proposal enables the interconnection of data and definitions across several domains utilizing the Semantic Web, articulated through RDF, Turtle, or JSON-LD languages. Tasks in decision-making difficulties are conventionally addressed using languages such as Python, R, or Scale, as well as more recently through cloud computing platforms. The portability of this code designed for a DM process is contingent upon dependencies related to libraries, development environments, or deployment rendering architectures. conversion alternative platforms problematic (see to Figure 1). Consequently, by utilizing services or service definitions in CC to compose a workflow instead of employing a programming language for migrating DM workflows, this issue is resolved: descriptions of these services are standardized, the of and intricacies execution are delegated to the CC platform. The prevailing trend in cloud computing is moving towards the industrialization of IT services, where standardization is essential to sustain the swift pace of change in IT and enhance the efficiency and effectiveness of service delivery. The primary goals of the article are the standardization and industrialization of DM services. The primary aim of this study is to delineate DM services for CC platforms while considering LD principles. This definition of the service encompasses not only the core components of DM (algorithms, workflows, parameters, or models) but also facilitates the definition and modeling of pricing, authentication, Service Level Agreements, computing resources (instances), and catalogues pertinent to the management of the CC introducedmccschema. service. We semantic framework based on Linked Data comprehensively define DM cloud facilitating services. their exchange, portability, search, and integration within computing. The document cloud organized as follows: This section presents the pertinent literature on the definition of cloud computing services. Section 3 presents

03779254 Page 194 of 203

our approach, the dmcc-schema, in full and illustrates all components along with their relationships. In section 4, we delineate many use cases for a real-world data mining service, offering a Random Forest algorithm as a service and addressing various elements pertinent to the service definition in cloud computing, including Service Level Agreements (SLAs) and pricing, along with a description of the algorithm. Section 5 delineates the conclusions and prospective endeavors.

II. RELATEDWORKS

Wang et al. (2020) proposed a semantic-based data mining framework in cloud environments to enhance service discovery and interoperability. Their model usedontologies to describe and match data mining services effectively.

Li and Zhang (2020) developed an ontology-driven approach for integrating heterogeneous data mining services in the cloud. Their approach improved the composition and reuse of services across different cloud platforms.

Kumar et al. (2021) introduced a semantic web service framework for data mining tasks in cloud computing. Their work focused on automating service selection and

improving service accuracy through semantic annotations.

Zhou et al. (2021) emphasized importance of semantic technologies for enhancing the scalability and flexibility of data mining in the cloud. Their work demonstrated that semantic frameworks could reduce service ambiguity and improve interoperability.Recent research trends highlight that the integration of semantic technologies, such as ontologies and semantic annotations, significantly improves the discovery, composition, and execution of mining services in cloud data computing.Despite advancements. challenges such as maintaining semantic consistency, handling dynamic service updates, and ensuring interoperability across heterogeneous cloud platforms remain active research areas.

III.SYSTEMANALYSIS

This significantly enriches the definition of the schema, allowing you to create the model definitions based on other existing schemata and vocabularies. LinkedUSDL, MEX, ML-Schema, OntoDM or Expose, can be considered when creating a workflow for DM service definition using LD. These proposals provide the definition of consistency in the main area of the service

03779254 Page 195 of 203

of Data Mining together with the LD properties, enabling the inclusion of other externals schemata that complement the key aspects of a CC service fully defined. With this review of state-of-the-art, a part of the range of services definition proposals has been studied to describe Cloud Computing services from different points of view.

PROPOSED SYSTEM

There are several proposals for the definition of services covering an important variety of both syntactic and semantic languages to achieve a correct definition and modeling of services. Solutions based on the proposal offered by Linked Data can solve the problem of defining services from a perspective more comprehensive. undertakes models and structures from the Semantic Web, a technology that aims to expose data on the web in a more reusable inter-operable way with and applications. The LD proposal allows you to link data and concept definitions from multiple domains through the use of the Semantic articulated with RDF languages. Tasks in DM problems are approached traditionally with languages such as Python, R, or Scale, among others and more recently from Cloud Computing platforms. The portability of this code created for an DM workflow is subject to

dependencies on libraries, development environments or even the deployment architecture, so migration to other platforms can be challenging Our scheme provides identity traceability to trace malicious cloud servers. If the cloud service providers exhibit any errors or illegal operations in the service process, users can trace back to the real identity of the corresponding cloud server based on the anonymous identity.

IV.IMPLEMENTATION

Modules:

Data Owners

Data owners are the owner and upload the files and pay for the resource consumption on file sharing. As the payers for cloud services, the data owners want the transparency of file utilization to ensure file security. The data owners require the cloud provider to justify the resource usage. In our system, the data owner is not always online.

Data Users

Data users want to obtain some files from the cloud provider stored on the cloud storage by following datamining services. They need to be authenticated by the cloud provider before the download (to thwart EDoS attacks). The authorized users then confirm (and sign for) the resource

03779254 Page 196 of 203

consumption for this download to the cloud provider.

Cloud Server

Cloud provider hosts the encrypted storage and is always online. It records the resource consumption and charges data owners based on that record. The cloud is not public-accessible in our system as it has an authentication based access control. Only data users satisfying the access policy can download the corresponding files. The cloud provider also collects the proof of the resource consumption to justify the billing.

Methodology

The proposed methodology aims to integrate semantic technologies with cloud-based data mining services to enhance service discovery, interoperability, and automation.

1.Requirement Analysis and Problem Identification

- Identify the challenges in discovering, integrating, and utilizing heterogeneous data mining services in cloud environments.
- Recognize the need for semantic descriptions to address issues such as service ambiguity.
- poor interoperability.
- Inefficient service composition.

2.Ontology Design for Data Mining Services

- Develop a domain-specific ontologyto formally describe data mining services, including:
- Service types (e.g., classification and clustering, regression)
- Input and output data formats
- Functional and non-functional attributes.
- Use standards such as OWL (Web
 Ontology Language) to model the
 ontology for better compatibility with
 existing semantic web technologies.

3. Semantic Annotation of Data Mining Services

- Annotate available cloud-based data mining services with semantic metadata based on the designed ontology.
- Define service properties, capabilities,input/output, specifications, and constraints using semantic annotations.
- Store service descriptions in a semantic repository or service registry for easy discovery.

03779254 Page 197 of 203

4.Semantic Service Discovery and Matchmaking

Implement a semantic reasoning engine to support intelligent service discovery.

- Enable users or applications to search for relevant data mining services by:
- Matching requested service requirements with semantic descriptions.
- Considering functional and non-functional parameters.
- Use ontology-based inference to resolve ambiguities
- Improve precision in service discovery.

5. Service Composition and Workflow Generation

- Allow dynamic composition of multiple data mining services to form complex workflows.
- Utilize semantic compatibility checks to ensure seamless integration between services.
- Automate workflow generation based on user requirements and service capabilities.

V.RESULTS AND DISCUSSION

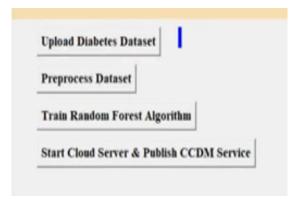


Fig 1 Home Page

This image shows the outlines of steps involved in developing and deploying a diabetes prediction service:

Data Management: Upload and preprocess a diabetes dataset.

Model Training: Train a machine learning model, specifically a Random Forest algorithm, on the prepared dataset.

Service Deployment: Start a cloud server and publish the developed model as a CCDM (Clinical Decision Support System) Service.



03779254 Page 198 of 203

Fig 2 Upload dataset Page

- **1.**The interface provides an option to upload files using the "Upload Files" button.
- **2.**It features a text area where information or logs can be displayed or entered.
- **3.**This layout is likely part of a data processing or machine learning application.
- **4.**Users can interact with the system by uploading datasets for further actions like training or prediction.

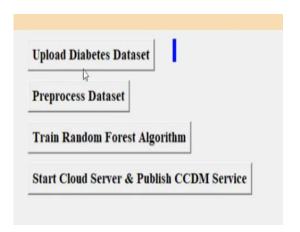


Fig 3 Preprocess Dataset Page

- 1. The interface includes a button labeled "upload files" for selecting and uploading data files.
- **2.** Below the button, there is a text display arealikely used for showing status messages or data previews.

- **3.** The layout is simple and user-friendly, designed for easy interaction with file-based workflows.
- **4.** This setup may be part of a machinelearning or data analysis application.
- **5.**Users can upload datasets here to begin processes such as preprocessing, model training, or prediction.

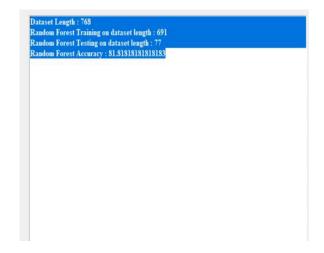


Fig 4 Preprocessed Pages

The image shows the output related to a machine learning model, likely a Random Forest Classifier, showing:

Dataset Length: The total size of the dataset used, which is 68.

Training Data: The Random Forest model was trained on a subset of the data with a length of 601 (this seems inconsistent with the total dataset length of 68, suggesting a

03779254 Page 199 of 203

potential typo or a different interpretation of "dataset length").

Testing Data: The model was tested on a dataset of length 7.

Accuracy: The Random Forest model achieved an accuracy of 81.81818181818183% during testing.

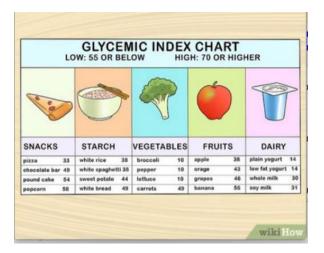


Fig 5 GI Chart Page

The image shows a Glycemic Index (GI) chart, which categorizes various foods based on their impact on blood sugar levels. The chart defines "Low GI" as 55 or below and "High GI" as 70 or higher.

The chart is organized into food categories:

Snacks: Includes items like pizza, chocolate bar, pound cake, and popcorn, with their respective GI values.

Starch: Features foods such as white rice, white spaghetti, sweet potato, and white bread, along with their GI values.

Vegetables: Lists vegetables like broccoli, pepper, lettuce, and carrots, with their corresponding GI values.

Fruits: Shows fruits including apple, orange, grapes, and banana, with their GI values.

Dairy: Presents dairy products like plain yogurt, low-fat yogurt, and whole milk, with their GI values.

This chart serves as a visual guide to help individuals make informed dietary choices by understanding how different foods may affect their blood glucose levels.

User Sense Data: 6,148,72,35,0,33.6,0.627,50 Abnormal Values. Predicted values: 1 Disease predicted as type 2 diabetes User Sense Data: 1.85.66.29.0.26.6.0.351.31 Normal Values. Predicted values: 0 No disease predicted User Sense Data: 8,183,64,0,0,23.3,0.672,32 Abnormal Values. Predicted values: 1 Disease predicted as type 2 diabetes User Seuse Data: 1,89,66,23,94,28.1,0.167,21 Normal Values. Predicted values: 0 No disease pre-User Sense Data: 1,189,60,23,846,30.1,0.398,59 Abnormal Values, Predicted values: 1 Disease predicted as type 2 diabetes Diet Plan Details Monday Breakfast: One poached egg and half a small avocado spread on one slice of Ezekiel bread, one orange. Total carbs: Lunch: Mexican bowl: two-thirds of a cup low-sodium canned pinto beans, 1 cup chopped spinach, a quarter cup chop bell peppers, 1 ounce (oz) cheese, 1 tablespoon (tbsp) salsa as sauce. Total carbs: Approximately 30. Snack: 20 1-gram baby carrots with 2 thsp hummus. Total carbs: Approximately 21. Dinner: 1 cup cooked lentil penne pasta, 1.5 cups veggie tomato sauce (cook garlic, mushrooms, greens, zucchini, and nd lean turkey. Total carbs: Approximately 35. Total carbs for the day: 125.

03779254 Page 200 of 203

Fig 6 Diet chart Page

This document appears to be a report related to health data analysis and diet planning, likely for an individual at risk of or with diabetes.

Disease Prediction: The report uses "User Sense Data" to predict the likelihood of "Disease" (specifically type 2 diabetes) and compares it to "Normal Values".

Abnormal Values: It identifies "Abnormal Values" in the user's data and predicts "1 Disease predicted as type 2 diabetes" based on those values.

Diet Plan Details: A detailed "Diet Plan" is provided for Monday, outlining specific meals (Breakfast, Lunch, Dinner) with estimated total carbohydrates for each, culminating in a "Total carbs for the Day: 125".

VI. CONCLUSION

This article presents dmcc-schema, a concise language for the description and definition of DM services in Cloud Computing. Our approach aims to consolidate, firstly, all elements pertaining to the definition of algorithms as a service for DM, and secondly, all other facets involved in the management of a Cloud Computing service, which are overlooked by alternative suggestions. Other service definition

approaches lack this integration; however, dmcc-schema addresses this deficiency by enabling the construction of comprehensive DM services for Cloud Computing settings. DMCC-schema introduced is lightweight instrument for modeling DM services, designed to provide a portable definition among various service providers. Consequently, the dmcc-schema enables the capture of all principal features and details (Cloud Computing administration and DM experimentation) of prominent CC providers such as Amazon, Azure, or Google. The dmcc-schema has been constructed based on the Semantic Web, utilizing an ontology language for its implementation and adhering to Linked Data principles for the reuse of other schemas, thereby enhancing the specified service modeling. Moreover, it guarantees that the definition of services can be enhanced and refined in the future, with the objective of providing a more adaptable definition ofservices the and accommodating changes in Cloud Computing management. The primary advantage of utilizing dmcc-schema is its ability to abstract several Data Mining service requirements from diverse Cloud Computing providers into a unified and generic specification, hence facilitating the standardization of Data Mining services.

03779254 Page 201 of 203

Consequently, the use and portability of services various across providers computing is guaranteed. Consequently, the distinctions among the definitions are harmonized, permitting the dmcc-schema to function, for instance, as a fundamental component of a Computing broker, which stores manages such Data Mining services from Cloud Computing providers. The scheme's usefulness is validated, since it enables the delineation of a DM service encompassing all its components (algorithms, costs/prices, etc.). The practical scenario utilizing the OC2DM deployment architecture provides hands-on capabilities for defining services with dmcc-schema, facilitating the composition and modeling DM workflows in Cloud Computing. The scheme's efficiency has been substantiated by transcribing actual DMCC services like Amazon SageMaker and confirming that the dmcc-schema encompasses all elements of these services, as evidenced by the CQs. With CQs, a series of questions are posed and answered to ascertain the scope of the problem, serving as a frequently utilized method for the validation of semantic frameworks. The validation component emphasizes both efficacy and efficiency. The practical situation demonstrates that

dmccschema significantly contributes to the standardization and industrialization of Data Mining services. Ultimately, as prospective extensions of the research presented in this will develop we semantic paper, specifications for Artificial Intelligence models, concentrating on Computational Intelligence, specifically evolutionary algorithms, neural networks, and fuzzy systems, to define them as services in Cloud Computing.

REFERENCES

[1] L. Liu, "Services computing: from cloud services, mobile services to internet of services," IEEE Transactions on Services Computing, vol. 9, no. 5, pp. 661–663, 2016.

[2] B. Marr, "Big data overload: Why most companies can't deal with the data explosion," Apr 2016. [Online]. Available: https://goo.gl/VZbe4R

[3] G. Inc., "Gartner says 6.4 billion connected 'things' will be in use in 2016," Gartner, Tech. Rep., 2016. [Online]. Available: https://www.gartner.com/newsroom/id/3165317

[4] T. K. Ho, "The random subspace method for constructing decision forests," IEEE transactions on pattern analysis and machine

03779254 Page 202 of 203

intelligence, vol. 20, no. 8, pp. 832-844, 1998.

[5] D. Lin, A. C. Squicciarini, V. N. Dondapati, and S. Sundareswaran, "A cloud brokerage architecture for efficient cloud service selection," IEEE Transactions on Services Computing, pp. 1–1, 2018.

[6] C. Bizer, T. Heath, and T. Berners-Lee, "Linked data-the story so far," International journal on semantic web and information systems, vol. 5, no. 3, pp. 1–22, 2009.

[7] T. Berners-Lee, J. Hendler, O. Lassila et al., "The semantic web," Scientific American, vol. 284, no. 5, pp. 28–37, 2001.

[8] G. Klyne and J. J. Carroll, "Resource description framework (rdf): Concepts and abstract syntax," W3C, 2006.

AUTHORS DETAILS



Mrs. J. KUMARI is an Assistant Professor in the Department of Master of Computer

Applications at QIS College of Engineering and Technology, Ongole, Andhra Pradesh. She earned Master of Computer Applications (MCA) from Osmania

University, Hyderabad, and her M.Tech in Computer Science and Engineering (CSE) from Jawaharlal Nehru Technological University, Kakinada (JNTUK). research interests include Machine Learning, programming language. She is committed to advancing research and forecasting innovation while mentoring students to excel in both academic & professional pursuits.



Mr.N.MADHUSUDHANA
RAO is a postgraduate
student pursuing a MCA in
the Department of

Computer Applications at QIS College of Engineering & Technology, Ongole an Autonomous college in Prakasam dist. He completed his undergraduate degree in B.Sc.(Tally, Multimedia) From (Acharya Nagarjuna University).

His academic interests include Cloud Computing, Artificial Intelligence, Cyber Security, and Data Structures.

03779254 Page 203 of 203