

# Semantics of Data Mining Services in Cloud Computing

Manuel Parra-Royon\*, Ghislain Ateazing<sup>†</sup> and Jose Manuel Benítez-Sanchez<sup>‡</sup>

January 15, 2019

## Abstract

In recent years with the rise of Cloud Computing (CC), many companies providing services in the cloud, are empowering a new series of services to their catalogue, such as data mining (DM) and data processing (DP), taking advantage of the vast computing resources available to them. Different service definition proposals have been put forward to address the problem of describing services in CC in a comprehensive way. Bearing in mind that each provider has its own definition of the logic of its services, and specifically of DM services, it should be pointed out that the possibility of describing services in a flexible way between providers is fundamental in order to maintain the usability and portability of this type of CC services. The use of semantic technologies based on the proposal offered by Linked Data (LD) for the definition of services, allows the design and modelling of DM services, achieving a high degree of interoperability. In this article a schema for the definition of DM services on CC is presented considering all key aspects of service in CC, such as prices, interfaces, Software Level Agreement (SLA), instances or DM workflow, among others. The new schema is based on LD, and it reuses other schemata obtaining a better and more complete definition of the services. In order to validate the completeness of the scheme, a series of DM services have been created where a set of algorithms such as *Random Forest* (RF) or *KMeans* are modeled as services. In addition, a dataset has been generated including the definition of the services of several actual CC DM providers, confirming the effectiveness of the schema.

## 1 Introduction

Cloud Computing has been introduced into our daily lives in a completely transparent and friction-less way. The ease of Internet access and the exponential increase in the number of connected devices has made it even more popular. Adopting the phenomenon of CC means a fundamental change in the way Information Technology (IT) services are explored, consumed or deployed. CC is a model of providing services to companies, entities and users, following the utility model, such as energy or gas. CC can be seen as a model of service provision where computer resources and computing power are contracted through the Internet of services (IS)[1]. A big part of the CC services providers are currently leveraging their wide computing infrastructure to offer a set of web services to enterprises, organizations and users.

The increase in the volume of data generated by companies and organizations is growing at an extremely high rate. According to Forbes [2], in 2020, the growth is expected to continue and data generation is predicted to increase by up to 4,300%, all motivated by the large amount of data generated by service users. By 2020, it is estimated that more than 25 billion devices will be connected to the Internet, according to *Gartner* [3], and that they will produce more than 44 billion GB of data annually. In this scenario, CC providers have successfully included data analysis services in their catalogue of services for massive processing of data and DM.

---

\*Department of Computer Sciences and Artificial Intelligence, Soft Computing and Intelligent Systems group, University of Granada, Spain.

<sup>†</sup>Department of R&D, Mondeca, Paris, France.

<sup>‡</sup>Department of Computer Sciences and Artificial Intelligence, Soft Computing and Intelligent Systems group, University of Granada, Spain.

These services allow the application of Artificial Intelligence (AI) and Machine Learning (ML) techniques on a large variety of data, offering an extensive catalogue of algorithms and workflows related to DM. Services, such as *Amazon SageMaker*<sup>1</sup> or *Microsoft Azure Machine Learning Studio*<sup>2</sup> (Table 1), offer a set of algorithms as services within CC platforms. Following this line, other CC platforms such as *Algorithmia*<sup>3</sup> or *Google Cloud ML*<sup>4</sup>, offer ML services at the highest level, providing specific services for the detection of objects in photographs and video, sentiment analysis, text mining or forecasting, for instance.

Each CC service provider offers a narrow definition of these services, which is generally incompatible with other service providers. For instance, where one provider has a service for RF algorithm, another provider has another name, features, or parameters for that algorithm, although the two might be the same. This makes it difficult to define services or service models independent of the provider as well as to compare services through a CC service broker [4]. Indeed, a standardization of the definition of services would boost competitiveness, allowing third parties to operate with these services in a totally transparent way, skipping the individual details of the providers. The effectiveness of CC would be greatly improved if there were a general standard for services definition [5].

There are several proposals for the definition of services. These proposals cover an important variety of both syntactic and semantic languages in order to achieve a correct definition and modelling of services. For the definition of this type of DM services, there is no specific proposal, due to the complexity of the services represented. Solutions based on the proposal offered by Linked Data [6] can solve the problem of defining services from a perspective more comprehensive. Linked Data undertakes models and structures from the Semantic Web, a technology that aims to expose data on the web in a more reusable and interoperable way with other applications. The Linked Data proposal allows you to link data and concept definitions from multiple domains through the use of the Semantic Web [7] articulated with *RDF* [8] or *Turtle* [9] languages.

The main objective of this work is the definition of DM services for CC platforms taking into account the Linked Data principles. The definition of the service is not only focused on the main part of the service (algorithms, workflow, parameters or models), but also allows the definition and modelling of prices, authentication, SLA, computing resources or catalogue. The *dmcc-schema* proposal provides a complete vocabulary for the exchange, consumption and negotiation of DM services in CC. With this schema it is possible to make queries in *SparQL* [10] about this type of CC DM services and obtain, for example, the set of providers that offer a certain algorithm, as well as the economic cost of the service. This allows for the comparison between different providers and DM services from varying points of view (costs, regions, instances, algorithms, etc). For the modelling of the different components of the service, existing schemata and vocabularies have been used, which have been adapted to the problem of the definition of DM services. In addition, new vocabularies have been created to cover specific elements of the domain that have not been available up until now. We present *dmcc-schema*, a proposal based on Linked Data to cover the entire definition of DM/ML cloud services and that allows the exchange, search and integration of this type of services in CC.

The paper is structured as follows: In the following section we present the related works in CC services definition. Section 3 put forward our proposal, the *dmcc-schema* in detail, and depicts all components with their interactions. We describe in section 4 some use cases for a real-world data mining service, providing with a RF algorithm as service and including some aspects related to the service definition in CC, such as SLA or prices of the service, and the algorithm description. Finally, we conclude our work in section 5.

---

<sup>1</sup><https://aws.amazon.com/sagemaker/>

<sup>2</sup><https://azure.microsoft.com/en-us/services/machine-learning-studio/>

<sup>3</sup><https://algorithmia.com/>

<sup>4</sup><https://cloud.google.com/products/machine-learning/>

## 2 Related work

DM and ML is a very highly relevant topic nowadays. These areas of knowledge provides entities, organizations and individuals with tools for data analysis and knowledge extraction. With the growth of CC and Edge Computing [11], DM services are taking a significant position in the catalogue offered by providers. The complexity of DM/ML services requires a complete specification, that is not limited to technical or conceptual considerations. The definition, therefore, must integrate key aspects of the service into the CC environment.

The definition of services in general has been approached from multiple perspectives. At the syntactic level with the reference of Services-Oriented-Architecture (SoA) [12] and XML (Extended Mark-Up Language) [13], derived languages such as WSDL [14], WADL [15], or SoAML [12], UDDI [16] (both related to UML [17]), it has been attempted to specify at the technical level the modelling of any Internet of Services.

The definition of services through semantic languages, enables to work with the modelling of services in a more flexible mode. Semantic languages allow to capture functional and non-functional features. OWL-S [18] integrates key factors for these services such as service discovery, process modelling, and service details. Web services can be described by using WSMO [19] as well. With WSMO it is possible to describe semantically the discovery, invocation and logic aspects of the service. Schemata such as SA-WSDL [20], which integrate both semantic and syntactic schema, associate semantic annotations to WSDL elements. SA-REST [21], WSMO-lite [22] or Minimal Service Model [23] (MSM), are also part of the group for the exchange, automation and composition of services, focused on the final service. On the other hand USDL [24] considers a global approach to service modelling that seeks to emphasize business, entity, price, legal, composition, and technical considerations. With Linked-USDL [25], this idea becomes concrete, in services integrated into the CC model. Linked-USDL assumes an important part of the key aspects that a service must provide, extending its usefulness to the ability to create CC services from scratch. With Linked-USDL, the modelling of entities, catalogue, SLA or interaction are considered.

The proposals address the modelling and definition of services in a generic mode, without dealing with the specific details of a DM/ML services in CC. For the treatment of these problems we typically use software platforms such as Weka [26], Knime [27] or *Orange* [28] and frameworks and libraries integrated into the most modern and efficient programming languages such as *C*, *Java*, *Python*, *R*, *Scala* and others. These environments lack elements for CC services modelling due to its nature.

There are several approaches to tackle an ontology-based languages for DM/ML services definition on experimentation and workflows. For example, *Exposé* [29] allows you to focus your work on the experimentation workflow [30]. With *OntoDM* [31] and using *BFO* (Basic Formal Ontology), it is aimed to create a framework for data mining experimentation and replication. In *MEXcore* [32] and *MEXalgo* [33] a complete specification of the process of description of experimentation in DM problems is made. In addition, together with *MEXperf* [32], which adds performance measures to the experimentation, the definition of this type of schema is completed. Another approach similar to *MEX* is *ML-Schema* [34], which attempts to establish a standard schema for algorithms, data and experimentation. *ML-Schema* consider terms like Task, Algorithm, Implementation, Hyper-Parameter, Data, Features or Models.

The most recent proposals try to unify different schemata and vocabularies previously developed following the Linked Data [7] guidelines. With Linked Data, you can reuse vocabularies, schemata and concepts. This significantly enriches the definition of the schema, allowing you to create the model definition based on other existing schemata and vocabularies. Linked-USDL, *MEX[core,algo,perf]*, *ML-Schema*, *OntoDM* or *Exposé* among others, can be considered when creating a workflow for DM/ML service definition using Linked Data. These proposals provide the definition of consistency in the main area of the service of DM together with the Linked Data 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 research material, a part of the range of services definition proposals have been studied in order to describe CC services from different points of view. In this proposal we seek

to bring together different proposals and some new ones in order to create a broad definition of DM services.

### 3 *dmcc-schema*: Data Mining services with Linked Data

Semantic Web applied to the definition of CC services, allow tasks such as negotiation, composition and invocation with a high degree of automation. This automation, based on Linked Data is fundamental in CC because it allows services to be discovered and explored for consumption by other entities using the full potential of *RDF* and *SparQL*. Linked Data [35] offers a growing body of reusable schemata and vocabularies for the definition of CC services of any kind [36].

In this article we propose *dmcc-schema*, a schema and a set of vocabularies which has been designed to address the problem of describing and defining DM/ML services in CC. Not only it focuses on solving the specific problem of modelling, with the definition of workflow and algorithms, but it also includes the main aspects of a CC service. *dmcc-schema* can be considered as a Linked Data proposal for DM/ML services. Existing Linked Data vocabularies have been integrated into *dmcc-schema* and new vocabularies have been created *ad-hoc* to cover certain aspects that are not implemented by other external schemata. Vocabularies have been re-used following Linked Data recommendations, filling important parts such as the definition of experiments and algorithms, as well as the interaction or authentication that were already defined in other vocabularies. A service offered from a CC provider, must have the following aspects for its complete specification:

- **Authentication.** The service or services require authenticated access.
- **Catalogue.** The provider has a catalogue of services ready to be discovered and used.
- **Entities.** Services interact between entities offering or consuming services.
- **Interaction.** Access points and interfaces for services consumption for users and entities.
- **Prices.** The services offered have a cost.
- **SLA/SLO.** The services have service level agreements and more issues.

There is no standardization about what elements a service in CC should have for its complete definition, but according to *NIST*<sup>5</sup>, it must meet aspects such as self-service (discovery), measurement (prices, SLA), among others. For the development of *dmcc-schema*, a comprehensive study of the features and services available on the DM/ML platforms on CC has been carried out. Table 1 contains information about the DM services analyzed from some CC providers; leading providers such as *Google*, *Amazon*, *Microsoft Azure*, *IBM* and *Algorithmia*, for whom, SLA, pricing for the different variants and conditions, service catalogue (DM/ML algorithms), methods of interaction with the service, and authentication have been studied.

Our schema has been complemented reusing external vocabularies: the interaction with the service, using *schema.org* [38], authentication, with *Web Api Authentication (waa)* [39], price design, with *GoodRelations* [40] and experimentation and algorithms with DM/ML, using the *dmcc-schema* vocabulary (*mls*). Table 3 shows the vocabularies reused in the *mls* schema to define each module of the full service. In table 2 a comparison between *dmcc-schema* and other schemata and vocabularies is performed.

The high level diagram of *dmcc-schema* entities can be seen in the figure 1. The detailed definition of each entity is developed in the following subsections.

*dmcc-schema* allows to structure the data and information of the services in a coherent way, normalizing the properties and the main concepts of the CC from the different DM service providers. Regarding

---

<sup>5</sup>National Institute of Standards and Technology, [www.nist.gov](http://www.nist.gov)

Table 1: CC Data Mining services analyzed.

Provider	Service name
<b>Google</b>	Cloud Machine Learning Engine
<b>Amazon</b>	Amazon SageMaker, Amazon Machine Learning
<b>IBM</b>	Watson Machine Learning, Data Science
<b>Microsoft Azure</b>	Machine Learning Studio
<b>Algorithmia</b>	Algorithms bundles

Table 2: Features included with compared with other.

	dmcc	mls[37]	Exposé[29]	MEX[32]	L-USDL[6]
Task	•	•	•	•	-
Impl.	•	•	•	-	-
Data	•	•	-	•	-
Model	•	•	•	•	•
Auth.	•	-	-	-	•
Pricing	•	-	-	-	•
SLA	•	-	-	-	•
Instan.	•	-	-	-	-

Table 3: Vocabularies reused by *dmcc-schema*.

Module	Reused Vocabularies
SLA	gr, schema, ccsla
Pricing	gr, schema, ccpricing, ccinstance, ccregion
Authentication	waa
Interaction	schema
Data Mining	mls, ccdm
Service Provider	gr, schema

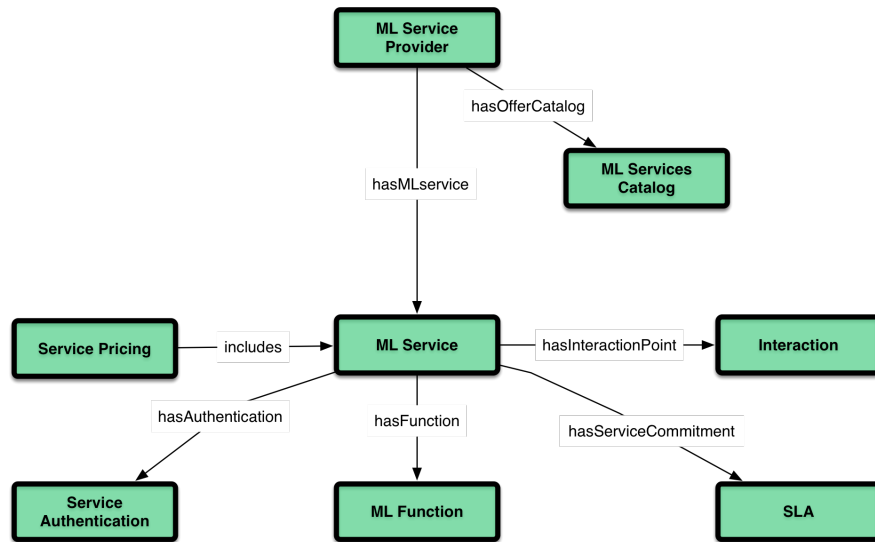


Figure 1: Main classes, and relations for *dmcc-schema* (*dmcc*).

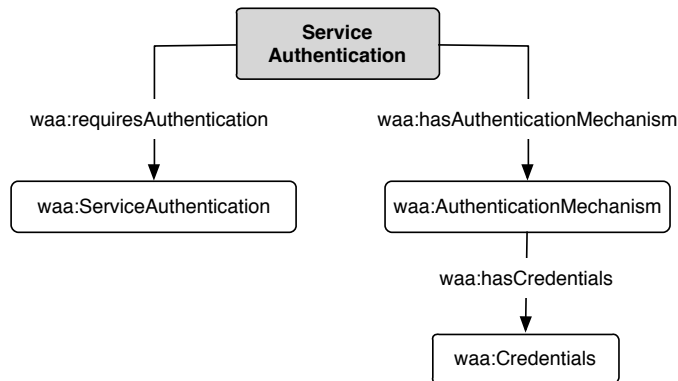


Figure 2: Authentication schema.

the benefits of using *dmcc-schema* over other schemata, a) offers a flexible data structure that can integrate the properties of existing DM services, b) unifies different schemata into one and offers friction less integration of the different schemata included, and c) comprises an all-in-one solution for the definition of CC services for DM. With *dmcc-schema* we have tried to emphasize the independence of the CC platform, that is, *dmcc-schema* does not define aspects of service deployment or elements closer to implementation tools. *dmcc-schema* definition is at the highest level and only addresses the definition and modelling aspects of the service without going into the details of infrastructure deployment. Below we detail each of the parts that compose our proposal.

### 3.1 Authentication

Nowadays, security in computer systems and in particular in CC platforms is an aspect that must be taken very seriously when developing CC services and applications. In this way, a service that is reliable and robust must also be secure against potential unauthorized access. In the outer layer of security in CC services consumption, authentication that should be considered as a fundamental part of a CC service definition.

Authentication on CC platforms covers a wide range of possibilities. It should be noted that for the vast majority of services, the most commonly used option for managing user access to services are *API Key* or *OAuth* and other mechanisms. *waa:WebApiAuthentication* has been included for authentication modelling. This schema allows to model many of the authentication systems available. In figure 2, the model for the definition of authentication services is depicted, together with the details of the authentication mechanisms.

### 3.2 Data Mining service

For the main part of the service, where the experimentation and execution of algorithms is specified and modelled, parts of *ML-Schema* (mls) have been reused. *MEXcore*, *OntoDM*, *DMOP* or *Exposé* also provide an adequate abstraction to model the service, but they are more complex and their vocabulary is more extensive. *ML-Schema* has been designed to simplify the modelling of DM/ML experiments and bring them into line with that which is offered by CC providers. We have extended *ML-Schema* by adapting its model to a specific one and inheriting all its features (figure 3).

The following vocabulary components are highlighted (*ccdm* is the name of the schema used):

- **ccdm:MLFunction** Set the operations, function or algorithm to be executed. For example **Random Forest** or **KNN**.
- **ccdm:MLServiceOutput** The output of the algorithm or the execution. Here the output of the experiment is modelled as **Model**, **Model Evaluation** or **Data**.

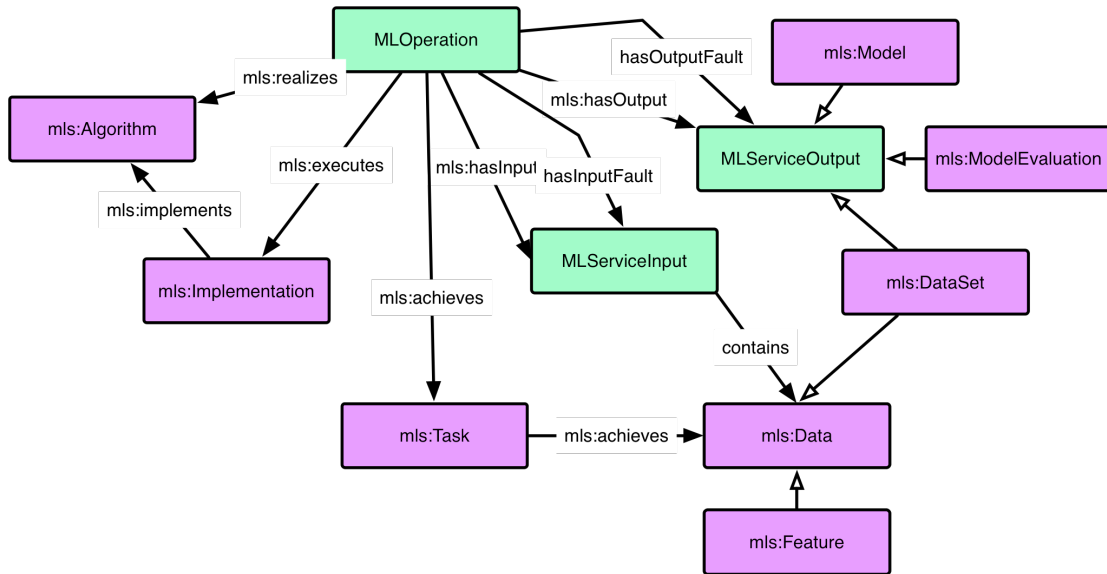


Figure 3: DM experimentation, workflow and algorithms execution.

- `ccdm:MLServiceInput` The algorithm input, which corresponds to the setting of the algorithm implementation. Here, you can describe the model the data entry of the experiment, such as the data set and parameters (`ccdm:MLServiceInputParameters`) of the algorithm executed.
- `mls:Model` Contains information specific to the model that has been generated from the run.
- `mls:ModelEvaluation` Provides the performance measurements of the model.
- `mls:Data` They contain the information of complete tables or only attributes (table column), only instances (row), or only a single value.
- `mls:Task` It is a part of the experiment that needs to be performed in the DM/ML process.

### 3.3 Interaction

The interaction with DM/ML services is generally done through a RESTful API. This API provides the basic functionality of interaction with the service consumer, who must be previously authenticated to use the services identified in this way. For the interaction the *Action* entity of the vocabulary *schema* was used. With this definition, the service entry points, methods and interaction variables are fully specified for all services specified by the API.

### 3.4 SLA: Software Level Agreement

The trading of CC services sets up a series of contractual agreements between the stakeholders involved in the services. Both the provider and the consumer of the service must agree on service terms. The service level agreements define technical criteria, relating to availability, response time or error recovery, among others. The SLO are specific measurable features of the SLA such as availability, throughput, frequency, response time, or quality of service. In addition with the SLO, it needs to contemplate actions when such agreements cannot be achieved where in this case compensation is offered.

The SLA studied for the DM/ML service environment are established by terms and definitions of the agreements ([41], [42], [43]). These terms may have certain conditions associated with them which, in the case of violation, involve a compensation for the guarantee service.

Table 4: Example of a MUP term and compensations related.

Monthly Uptime Percentage (MUP)	Compensation
[0.00%, 99.00%]	25%
[99.00%, 99.99%]	10%

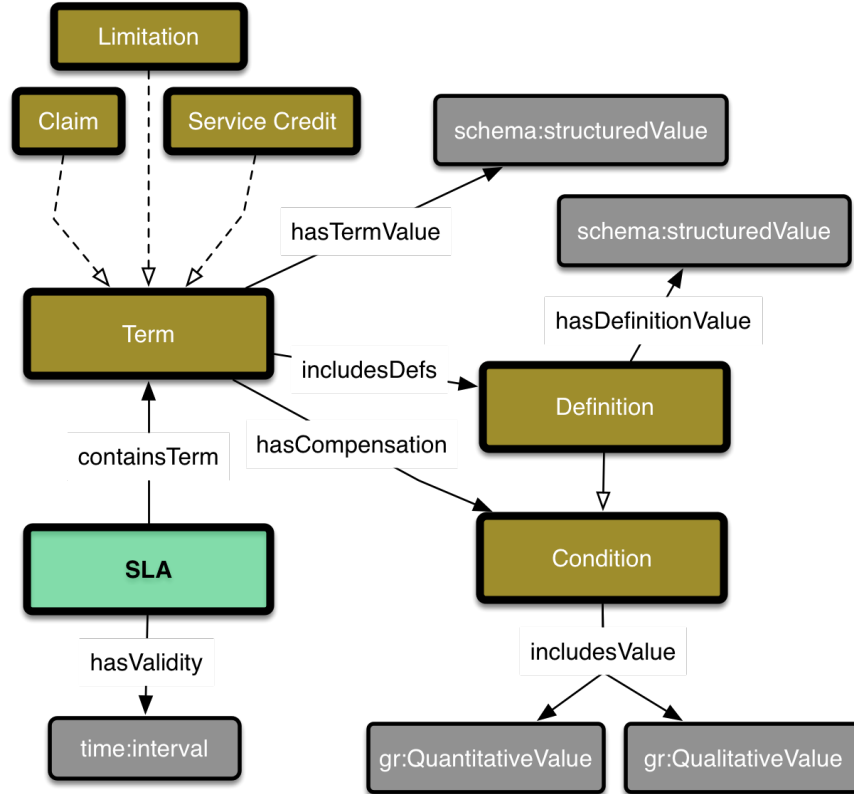


Figure 4: SLA schema for CC services (*ccsla*).

In general terms, SLA for CC services is given by a term of the agreements that contains one or more definitions similar to *Monthly Up-time Percentage* (MUP), which specifies the maximum available minutes less downtime divided by maximum available minutes in a billing month. For this term, a metric or interval is established over which a compensation is applied in the case that the agreement term is not satisfied. In this context a schema named *ccsla*<sup>6</sup> has been created for the definition of all the components of an SLA. In Table 4 an example with a pair of MUP intervals and compensation related is shown.

In figure 4 the complete schema of the SLA modelling that has been designed is illustrated.

### 3.5 Pricing

Like any other utility-oriented service, CC services are affected by costs and pricing that vary with the price of the service offered. Following the pay-as-you-go model, the costs of using the service are directly related to the characteristics of the service and the use made of it. The pricing of services in CC is a complex task. Not only there is a price plan for temporary use of resources, but it is also affected by technical aspects, configuration or location of the service. Table 5 collects some of the price modifying

<sup>6</sup>Available at <http://lov.okfn.org/dataset/lov/vocabs/ccsla>



Table 5: Price components for each CC provider analyzed.

Provider / Service	Region	Instances	Storage
Amazon SageMaker	y	y	y
Amazon Lambda	y	n	n
Azure Machine Learning	y	y	y
Azure Functions	y	n	n
IBM BlueMix	y	y	n
Google Machine Learning	y	y	n
Algorithmia	n	n	n

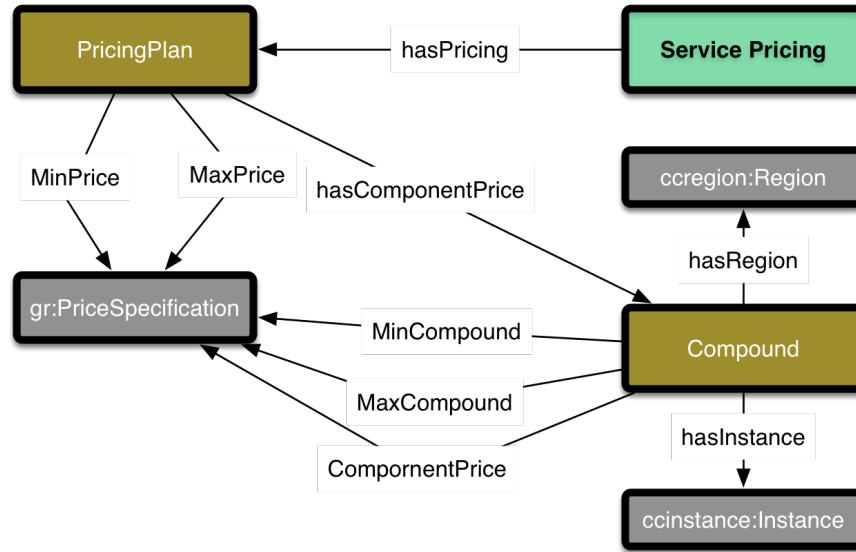


Figure 5: Pricing schema with *ccpricing*.

aspects of the use of the DM/ML service. In the specific segment of DM services in CC there are several elements to consider.

The following components of the vocabulary stand out from the diagram of figure 5:

- **ccpricing:PricingPlan**. It allows to define the attributes of the price plan and their details. For instance, a free, premium plan, or similar. It must match the attributes of *MinPrice* and *MaxPrice*.
- **gr:PriceSpecification**. It provides all the tools to define the prices of the services at the highest level of detail. Comes from *GoodRelations* (*gr*) schema.
- **ccpricing:Compound**. Components of the price of the service. Services are charged according to values that are related to certain attributes of the service configuration. In compound you can add aspects such as type of instances, cost of the region, among others.
- **ccinstances:Instance**. Provides the specific vocabulary for the definition of price attributes related to the instance or instances where the DM service is executed.
- **ccregions:Region**. It allows to use of locations and regions schema, offered by providers. This is an additional value to the costs of the service, as part of the price.

For the modelling of this part, three vocabularies supplementing the definition of prices, have been developed which were not available. On the one side *ccinstances*<sup>7</sup>, which contains everything necessary for the definition of instances (*CPU*, number of *CPU* cores, model, *RAM* or *HD*, among others), on the other side *ccregions*<sup>8</sup>, which provides the general modelling of the service regions and location in CC and the last one the schema for the prices modelling *ccpricing*<sup>9</sup>.

The complete diagram of the *dmcc-schema* is available on the project website [44], from which you can explore the relationships, entities, classes, and attributes of each of the modules that comprise the proposal.

## 4 Use cases and validation

For reasons of space we will not detail all the data or attributes in depth and will only consider what is most important for the basic specification of the service and his comprehension. All detailed information on each component, examples and data sets are available on the project's complementary information website.

In this section we will illustrate how *dmcc-schema* can be used effectively to define a specific DM service, as well as additional aspects related to CC such as SLA, interaction, pricing or authentication. To do this we will create a service instance for a RF [45] algorithm for classification. In RF service implementation, we will use the default function model and parameters of this algorithm in R [46] language. We have chosen the parameters format of the R algorithms, due to the great growth that this programming language is currently having for data science. In order to instantiating of the service we will use vocabularies such *dcterms* (*dc:*) [47], *GoodRelations* (*gr:*) or *schema.org* (*s:*). In these cases of use we will name *dmcc-schema* as *dmcc* namespace.

The entry point for defining the service is the instantiating of the *dmcc:ServiceProvider* entity, including its data, using the class *dmcc:MLServiceProvider*. We will specify an example of DM service hosted on our CC platform *dicits.ugr.es*.

```

1  _:MLProvider a dmcc:MLServiceProvider;
2  rdfs:label "ML Provider"@en ;
3  dc:description
4     "DICITS ML SP"@en ;
5  gr:name "DITICS ML Provider";
6  gr:legalName "U. of Granada";
7  gr:hasNAICS "541519";
8  s:url <\protect\vrule width0pt\protect\href{http://www.dicits.ugr.es}{http://www.dici
9  s:serviceLocation
10 [ a s:PostalAddress;
11   s:addressCountry "ES";
12   s:addressLocality "Granada";
13 ] ;
14 s:contactPoint
15 [
16   a s:ContactPoint;
17   s:contactType "Costumer Service";
18   s:availableLanguage [
19     a s:Language;
20     s:name "English";];
21   s:email "ml@dicits.ugr.es";
22 ] ;
23 dmcc:hasMLService

```

<sup>7</sup><http://lov.okfn.org/dataset/lov/vocabs/cci>

<sup>8</sup><http://lov.okfn.org/dataset/lov/vocabs/ccr>

<sup>9</sup><http://lov.okfn.org/dataset/lov/vocabs/ccp>

```

24         _:MLServiceDicitsRF ,
25         _:MLServiceDicitsKMeans ;
26 dmcc:hasOfferCatalog
27         _:MLServiceDicitsCatalog ;
28 .

```

Listing 1: New DM service definition.

The code above shows the specific details of the service provider, such as legal aspects, location or contact information. The service provider `dmcc:MLServiceProvider` has the property `dmcc:hasMLService`. With this property two services such as `_:MLServiceDicitsRF` and `_:MLServiceDicitsKMeans` (listing 1 line 24-25) are defined.

Listing 2 shows the instantiating of each of the entities for `_:MLServiceDicitsRF`. `_:MLServiceDicitsRF` allows defining all the service entities such as SLA (listing 2 line 8), algorithm (listing 2 line 10), prices (listing 2 line 15), interaction points (listing 2 line 6) and authentication (listing 2 line 12).

```

1  _:MLServiceDicitsRF a dmcc:MLService ;
2  rdfs:label
3    "ML Service dicits.ugr.es"@en ;
4  dc:description
5    "DICITS ML Service"@en ;
6  dmcc:hasInteractionPoint
7    _:MLServiceInteraction ;
8  dmcc:hasServiceCommitment
9    _:MLServiceSLA ;
10 dmcc:hasFunction
11  _:MLServiceFunction ;
12 dmcc:hasAuthentication
13  _:MLServiceAuth ;
14 dmcc:hasPricingPlan
15  _:MLServicePricing ;
16 .

```

Listing 2: Components of the DM service

The interaction with the service is performed using the `dmcc:Interaction` class that includes the property `dmcc:hasEntryPoint` that allows to define an `Action` on a resource or object. In this case we use the vocabulary `schema.org` to set the method of access to the RF service, for which we specify that the service will be consumed through a RESTful API, with HTTP POST access method, the URL template `http://dicits.ugr.es/ml/rf/`, for instance and its parameters.

The selected service requires authentication. An API KEY authentication has been chosen. The instantiating of service authentication is defined requesting authentication with `waa:requiresAuthentication` and setting up the authentication mechanism `waa:hasAuthenticationMechanism` to `waa:Direct` and as credentials an API KEY is also used as shown in listing 3.

```

1  _:MLServiceAuth
2    a dmcc:ServiceAuthentication ;
3  rdfs:label
4    "Service Authentication"@en ;
5  dc:description "Service Auth"@en ;
6  waa:requiresAuthentication waa:All ;
7  waa:hasAuthenticationMechanism
8    [ a waa:Direct ;
9      waa:hasInputCredentials
10     [ a waa:APIkey ;
11       waa:isGroundedIn "key" ;
12     ] ;

```

```

13   waa:wayOfSendingInformation
14       waa:ViaURI;
15 ]
16 .

```

Listing 3: Service authentication.

For the definition of the SLA, in our example service we have taken the some providers as Amazon or Azure where they identify a term Month Uptime Percentage (MUP). For the cases where this occurs SLA define two ranges: less than *99.00* is compensated by *30* credits (in service usage) and the range of *99.00* to *99.99* is compensated with *10* credits. To model the SLA we use `ccsla:SLA` together with the property `ccsla:cointainsTerm` `_:SLATermMUP_A`; in which we define the specific terms.

To define the range of the term is used in `_:SLADefinition_A`, as shown in listing 4.

```

1  _:SLADefinition_A a ccsla:Definition;
2  ccsla:hasDefinitionValue [
3    a s:structuredValue;
4    s:value [
5      a s:QuantitativeValue;
6      s:maxValue 99.99;
7      s:minValue 99.00;
8      s:unitText "Percentaje";
9    ];
10 ];
11 .

```

Listing 4: SLA term definition.

The schema for SLA (*ccsla*) and other examples with the SLA for *Amazon Web Services* and *Microsoft Azure* can be accessed from the website of the project [44].

For the definition of the economic cost of the service we have considered two variants for the example. The first is for free service use, limited to *250* hours of execution of algorithms within an instance (Virtual Machine) with a CPU model *Intel i7*, 64 GB of RAM and one region. The second pricing model where you charge what you consume for the service in *USD/h.*, for an instance and one region.

We can define multiple pricing plans, for this example a free plan is specified with `_:ComponentsPricePlanFree`. The price modelling is done with our proposal using the definition of prices provided by *ccpricing*.

For each price plan we take into account the variables and features that affect the price. These are: region, instance type and other components using `_:ComponentsPricePlanFree`.

For example, to define features of the type of instance used in the free plan, we use `ccinstance:Instance`; and a few attributes like RAM or CPU as seen in listing 5. More examples for the schema *ccinstances*, including a small dataset of Amazon instances [44] (types *T1* and *M5*) are available in the web site of the complementary information of the paper.

```

1  _:InstanceFree
2    a ccinstances:Instance;
3    ccinstances:hasRAM [
4      a ccinstances:ram;
5        s:value "64"
6        s:unitCode "E34";
7    ] ;
8    ccinstances:hasCPU [
9      a ccinstances:cpu;
10     ccinstances:cpu_model
11     "Intel i7";
12   ] ;
13 .

```

Listing 5: Pricing and instances.

In order to define the `_:MaxUsageFree`, we need to determine the free access plan to the service and the limitation of compute hours to 250. For this purpose we use `gr:PriceSpecification` and `gr:Offering` classes as shown in listing 6.

```

1  _:MaxUsageFree a
2    gr:PriceSpecification ,
3    gr:Offering;
4  gr:max 0.00;
5  gr:priceCurrency "USD";
6  gr:includesObject [
7    a gr:TypeAndQualityNode;
8    gr:amountOfThisGood "250";
9    gr:hasUnitOfMeasurement "HRS";
10 ];
11 .

```

Listing 6: Price specification for the free plan.

Additional examples for the *ccpricing* schema, and a dataset of *Amazon SageMaker* pricing plans, are available in [44].

Finally, we define the service algorithm using `ccdm:MLFunction` for the definition of a RF function, where we specify the input parameters (data set and hyper-parameters), the output of the algorithm, among others (see listing 7).

```

1  _:RandomForest_Function
2    a ccdm:MLFunction ;
3    ccdm:hasInputParameters
4      _:RF_InputParameters;
5    mls:hasInput
6      _:RF_Input;
7    mls:hasOutput
8      _:RF_Output .

```

Listing 7: Operations for the algorithm/function.

The input and output data of the algorithms must be included in the definition of the data mining operation to be performed. The input of data, which can be parameters `_:RF_InputParameters` or data sets `_:RF_Input`. Input parameters of the algorithm can be defined with `dmc:MLServiceInputParameters` and the parameter list `_:parameter_01`, [...].

Definition of `ccdm:hasInputParameters` `_:RF_InputParameters` allows you to specify the general input parameters of the algorithm. For example for RF `dc:title "ntrees"` (number of trees generated), as well as whether `ccdm:mandatory "false"` is mandatory and its default value, if it exists. Listing 8 shows the definition of one of the parameters `parameter_01`. The other algorithm parameters are defined in the same way.

```

1  _:parameter_01
2    a ccdm:MLServiceInputParameter ;
3    ccdm:defaultvalue "100" ;
4    ccdm:mandatory "false" ;
5    dc:description
6      "Number of trees" ;
7    dc:title "ntrees" .

```

Listing 8: Example of parameter and features.

An `mls:Model` model and an evaluation of the `mls:ModelEvaluation` model have been considered for specifying the results of the RF service execution in `mls:hasOutput` `_:RF_Output`. Model evaluation is the specific results if the algorithm returns a value or set of values. When the service algorithm is

pre-processing the result is a data set. For the model you have to define for instance whether the results are PMML [48] with `_:RF_Model a dmc:PMML_Model` as shown in listing 9.

```
1 _:KMeans_Model
2   a ccdm:PMML_Model ;
3     ccdm:storagebucket
4       <dicits://models/> ;
5     dc:description
6       "PMML model" ;
7     dc:title "PMML Model" .
```

Listing 9: PMML Model and storage of the service output.

Other services implemented on *dmcc-schema* as examples, such as *NaiveBayes*, *LinearRegression*, *SVM* or *Optics*, are available in the supplementary website of the paper [44].

In order to validate the scheme, a dataset has been created containing the *dmcc-schema* description of several CC providers and their respective algorithms as services. For each provider, the specific data of the service has also been described, such as regions, instances and economic cost of each one of the variants of the consumption of the services. To confirm the validity of *dmcc-schema*, multiple queries with *SparQL* have been carried out to extract information from the dataset, such as to know the providers that offer a specific DM service or to obtain the best price to run a RF algorithm. Dataset and query results can be checked from the project site [44].

## 5 Conclusion

In this article we have presented *dmcc-schema*, a simple and direct schema for the description and definition of DM services in CC. Our proposal tries to gather, on the one hand, everything related to the definition of the experimentation, workflow and algorithms and on the other hand, all the other aspects that compose a complete CC service. Our schema has been built on the basis of Semantic Web, using an ontology language to implement it and following the Linked Data directives regarding the re-use of other schemata, which perfectly enrich the service modelling that has been designed.

*dmcc-schema* is presented as a light-weight tool for services modelling that allows the creation of a complete DM service that includes all providers of the CC platforms, adapting in a flexible way to the differences of definition and description of services of the most well-known providers.

The example of use shown, illustrates the effortless definition of a service whose objective is to execute a simple RF algorithm, and indicating other aspects related to the CC service itself.

One of the advantages of using *dmcc-schema* is that it abstracts differences between heterogeneous CC providers for DM services in order to have a single and unique specification that can bring together different services specifications. In this way the differences between the definitions are balanced, allowing to *dmcc-schema* to be used as an integral part of a CC Services Broker, storing such services from different providers.

Finally, it is important to highlight that *dmcc-schema* is being used successfully within a computing and workflow platform for DM, called OCCML. As part of the platform, *dmcc-schema* is used to define and describe complete DM services, allowing a high degree of flexibility and portability.

## Acknowledgment

M. Parra-Royon holds a "Excelencia" scholarship from the Regional Government of Andalucía (Spain). This work was supported by the Research Projects *P12-TIC-2958* and *TIN2016-81113-R* (Ministry of Economy, Industry and Competitiveness - Government of Spain).

## 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://www.forbes.com/sites/bernardmarr/2016/04/28/big-data-overload-most-companies-cant-deal-with-the-data-explosion/#4f478c0e6b0d>
- [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] 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.
- [5] S. Ghazouani and Y. Slimani, “A survey on cloud service description,” *Journal of Network and Computer Applications*, vol. 91, pp. 61–74, 2017.
- [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.
- [9] D. Beckett, “Turtle-terse rdf triple language,” <http://www.ibr.t.bris.ac.uk/discovery/2004/01/turtle/>, 2008.
- [10] E. Prud, A. Seaborne *et al.*, “Sparql query language for rdf,” *W3C*, 2006.
- [11] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, “Edge computing: Vision and challenges,” *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637–646, 2016.
- [12] E. Newcomer and G. Lomow, *Understanding SOA with Web services*. Addison-Wesley, 2005.
- [13] T. Bray, J. Paoli, C. M. Sperberg-McQueen, E. Maler, and F. Yergeau, “Extensible markup language (xml).” *World Wide Web Journal*, vol. 2, no. 4, pp. 27–66, 1997.
- [14] E. Christensen, F. Curbera, G. Meredith, and S. Weerawarana, “Web Service Definition Language (WSDL),” World Wide Web Consortium, Tech. Rep., Mar. 2001. [Online]. Available: <http://www.w3.org/TR/wsdl>
- [15] M. J. Hadley, “Web application description language (wadl),” *W3C*, vol. 23, 2006.
- [16] T. Bellwood, L. Clément, D. Ehnebuske, A. Hately, M. Hondo, Y. L. Husband, K. Januszewski, S. Lee, B. McKee, J. Munter *et al.*, “Uddi version 3.0,” *Published specification, Oasis*, vol. 5, pp. 16–18, 2002.
- [17] A. Dennis, B. H. Wixom, and D. Tegarden, *Systems analysis and design: An object-oriented approach with UML*. John wiley & sons, 2015.
- [18] D. Martin, M. Burstein, J. Hobbs, O. Lassila, D. McDermott, S. McIlraith, S. Narayanan, M. Paolucci, B. Parsia, T. Payne *et al.*, “Owl-s: Semantic markup for web services,” *W3C member submission*, vol. 22, pp. 2007–04, 2004.

- [19] J. Domingue, D. Roman, and M. Stollberg, “Web service modeling ontology (wsmo)-an ontology for semantic web services,” 2005.
- [20] J. Kopecký, T. Vitvar, C. Bournez, and J. Farrell, “SawSDL: Semantic annotations for WSDL and XML schema,” *IEEE Internet Computing*, vol. 11, no. 6, 2007.
- [21] M. Klusch, “Service discovery,” in *Encyclopedia of Social Network Analysis and Mining*. Springer, 2014, pp. 1707–1717.
- [22] J. Kopecký, T. Vitvar, D. Fensel, and K. Gomadam, “hRESTS & microwSMO,” *STI International, Tech. Rep.*, 2009.
- [23] M. Taheriyani, C. A. Knoblock, P. Szekely, and J. L. Ambite, “Rapidly integrating services into the linked data cloud,” in *International Semantic Web Conference*. Springer, 2012, pp. 559–574.
- [24] S. Kona, A. Bansal, L. Simon, A. Mallya, G. Gupta, and T. D. Hite, “USDL: A service-semantics description language for automatic service discovery and composition1,” *International Journal of Web Services Research*, vol. 6, no. 1, p. 20, 2009.
- [25] C. Pedrinaci, J. Cardoso, and T. Leidig, “Linked USDL: a vocabulary for web-scale service trading,” in *European Semantic Web Conference*. Springer, 2014, pp. 68–82.
- [26] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, “The Weka data mining software: An update,” *SIGKDD Explor. Newsl.*, vol. 11, no. 1, pp. 10–18, Nov. 2009. [Online]. Available: <http://doi.acm.org/10.1145/1656274.1656278>
- [27] M. R. Berthold, N. Cebron, F. Dill, T. R. Gabriel, T. Kötter, T. Meinl, P. Ohl, C. Sieb, K. Thiel, and B. Wiswedel, “KNIME: The Konstanz Information Miner,” in *Studies in Classification, Data Analysis, and Knowledge Organization (GfKL 2007)*. Springer, 2007.
- [28] J. Demšar, T. Curk, A. Erjavec, Črt Gorup, T. Hočevar, M. Milutinovič, M. Možina, M. Polajnar, M. Toplak, A. Starič, M. Štajdohar, L. Umek, L. Žagar, J. Žbontar, M. Žitnik, and B. Zupan, “Orange: Data mining toolbox in python,” *Journal of Machine Learning Research*, vol. 14, pp. 2349–2353, 2013. [Online]. Available: <http://jmlr.org/papers/v14/demsar13a.html>
- [29] J. Vanschoren and L. Soldatova, “Exposé: An ontology for data mining experiments,” in *International workshop on third generation data mining: Towards service-oriented knowledge discovery (SoKD-2010)*, 2010, pp. 31–46.
- [30] F. Marozzo, D. Talia, and P. Trunfio, “A workflow management system for scalable data mining on clouds,” *IEEE Transactions on Services Computing*, 2016.
- [31] P. Panov, S. Džeroski, and L. Soldatova, “OntoDM: An ontology of data mining,” in *Data Mining Workshops, 2008. ICDMW’08. IEEE International Conference on*. IEEE, 2008, pp. 752–760.
- [32] D. Esteves, D. Moussallem, C. B. Neto, T. Soru, R. Usbeck, M. Ackermann, and J. Lehmann, “Mex vocabulary: a lightweight interchange format for machine learning experiments,” in *Proceedings of the 11th International Conference on Semantic Systems*. ACM, 2015, pp. 169–176.
- [33] D. Esteves, D. Moussallem, C. B. Neto, J. Lehmann, M. C. R. Cavalcanti, and J. C. Duarte, “Interoperable machine learning metadata using mex.” in *International Semantic Web Conference (Posters & Demos)*, 2015.
- [34] D. Esteves, A. Lawrynowicz, P. Panov, L. Soldatova, T. Soru, and J. Vanschoren, “ML Schema Core Specification,” World Wide Web Consortium, Tech. Rep., Oct. 2016. [Online]. Available: <http://ml-schema.github.io/documentation/ML%20Schema.html>



- [35] T. Heath and C. Bizer, “Linked data: Evolving the web into a global data space,” *Synthesis lectures on the semantic web: theory and technology*, vol. 1, no. 1, pp. 1–136, 2011.
- [36] P.-Y. Vandenbussche, G. A. Atemezing, M. Poveda-Villalón, and B. Vatant, “Linked open vocabularies (lov): a gateway to reusable semantic vocabularies on the web,” *Semantic Web*, vol. 8, no. 3, pp. 437–452, 2017.
- [37] “Mlschema: machine learning and data mining ontologies,” <http://ml-schema.github.io/documentation/ML%20Schema.html#mapping>, accessed: 2018-05-20.
- [38] R. V. Guha, D. Brickley, and S. MacBeth, “Schema.org: Evolution of structured data on the web,” *Queue*, vol. 13, no. 9, pp. 10:10–10:37, Nov. 2015. [Online]. Available: <http://doi.acm.org/10.1145/2857274.2857276>
- [39] M. Maleshkova, C. Pedrinaci, J. Domingue, G. Alvaro, and I. Martinez, “Using semantics for automating the authentication of web apis,” *The Semantic Web–ISWC 2010*, pp. 534–549, 2010.
- [40] M. Hepp, “Goodrelations: An ontology for describing products and services offers on the web,” *Knowledge Engineering: Practice and Patterns*, pp. 329–346, 2008.
- [41] B. Ambulkar and V. Borkar, “Data mining in cloud computing,” in *MPGI National Multi Conference*, vol. 2012, 2012.
- [42] X. Li and J. Du, “Adaptive and attribute-based trust model for service-level agreement guarantee in cloud computing,” *IET Information Security*, vol. 7, no. 1, pp. 39–50, 2013.
- [43] Z. Zheng, J. Zhu, and M. R. Lyu, “Service-generated big data and big data-as-a-service: an overview,” in *IEEE International Congress on Big Data (BigData Congress)*. IEEE, 2013, pp. 403–410.
- [44] M. Parra, “Data mining service definition in cloud computing.” [Online]. Available: <http://dicits.ugr.es/linkedata/dmservices/>
- [45] T. K. Ho, “The random subspace method for constructing decision forests,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 20, no. 8, pp. 832–844, 1998.
- [46] R. C. Team *et al.*, *R: A language and environment for statistical computing*, 2013.
- [47] S. Weibel, J. Kunze, C. Lagoze, and M. Wolf, “Dublin core metadata for resource discovery,” Network Working Group, Tech. Rep., 1998.
- [48] A. Guazzelli, M. Zeller, W.-C. Lin, G. Williams *et al.*, “Pmml: An open standard for sharing models,” *The R Journal*, vol. 1, no. 1, pp. 60–65, 2009.