

# FROM SEMANTIC WEB TO FACETED CLASSIFICATION

## A CASE STUDY AND FIVE LINES OF FUTURE RESEARCH

**Carlo Bianchini**

*Department of Musicology and Cultural  
Heritage, University of Pavia, Pavia, Italy*

### **KEYWORDS:**

*Faceted classification, Semantic web, Colon Classification, data.bnf.fr, Italian Literature classification.*

### **ABSTRACT**

*A case study about the automatic creation of faceted class numbers for special topics starting from data available on the semantic web is described. The aim is to investigate two questions: 1) Is the automatic, quick, cost-effective, large-scale production of a Schedule of Colon Classification class numbers for literary classics possible?, and 2) More generally: which are the requirements for an automatic production of class numbers of a faceted scheme of classification? What kind of research is necessary to develop such a process?*

*The case study focuses on data extracted from data.bnf.fr and on their reuse to get correct and complete Colon Classification class numbers for a sample of Italian literature writers (born from 1800 to 1900). The process of creation of class numbers starting from data identification and extraction, through their cleanup, transformation and translation in classified notation by means of Open Refine is presented. Case study results and their relevance for starting questions are discussed and research questions for possible future developments are identified.*

## Introduction

Ranganathan's desire to overcome existing classification schemes and to create a new one based on new semantic and syntactical features is the foundational idea of *Colon Classification* (CC) (Ranganathan, 1963). In the new classification scheme, simple and composed subjects<sup>1</sup> are to be expressed by means of semantic and syntactical tools represented in the facet formula. Ranganathan states this basic idea by the contraposition of schemes *of* classification (e.g. Enumerative Classifications) – and schemes *for* classification (e.g. Faceted Classifications) (Ranganathan, 1965, pp. 13-14).

Enumerative Classifications are founded on an articulated scheme that contains all the necessary numbers for the most parts of the subjects. Ranganathan explains that

An Enumerative Scheme for classification consists essentially of a single schedule enumerating all subjects – of the past, the present and the anticipatable future.

Such a schedule will have necessarily to be long. Further, it will soon be overpowered by the emergence of subjects beyond the anticipated ones – overpowered in the sense that it would be difficult to place each new subject in its filial position among the already existing subjects (Ranganathan, 1967).

In this sense, *Library of Congress Classification* (LCC) and Rider's *International Classification* are designed as Enumerative Classifications. *Dewey Decimal Classification* (DDC) is an Almost-Enumerative Scheme for Classification, as it has in addition a few schedules of common isolates (*Tables*). Stefano Tartaglia underlines that

Enumeration implies either a deliberate limitation of the number of the classes and, in this case, it is functional just for a non-analytical classification, or the provision of every combination of concepts that might constitute the subject of a document; but the latter is impossible,

1 A simple subject is a subject constituted by one 'isolate' – that is one idea, that can be represented by one authorized term and one notation composed by one or more figures. For example, the subject of a book about *mathematics*, that is represented by the idea of Mathematics, expressed by the term *mathematics* and, in CC, by the figure "B" is a simple (or basic) subject. Instead, composed subjects are constituted by two or more isolates, and, depending on the kind of relationships among these isolates, can be compound or complex subject. For example, "History of mathematics in Europe during the XVII Century" and "Encyclopedia of Italian linguistics" are two composed subjects; on the contrary, "Statistics for librarians" and "Influences of Dante on Leopardi" are two complex (or phased) subjects (Ranganathan, 1959, 1967).

so that the provision results in a selection – always partially aleatory – of the most probable subjects (Tartaglia, 1998, p. 12).<sup>2</sup>

CC is an analytic-synthetic classification, and it does not supply ready-made classification numbers, but *Standard schedules* that provide a list of ‘isolates’ (i.e., the base elements of the ‘Meccano’ that constitutes the classification) to be used to compose a class number representing the simple or composed *specific subject* of a resource.<sup>3</sup> So, if the former type of classifications is based on an aleatory anticipated prevision of any possible subject to be classified (and in this sense they are scheme of classification), faceted classifications provide just the basic tools to compose each time a number containing all the necessary isolates (so that they are a scheme for classification). By a different process of synthesis based on the numbers listed in the schedules, class numbers for any simple or composed subject can be built by combinations and permutations allowed by syntax rules (such as the Facet Formula placed at the beginning of each Main Class Schedule). Ranganathan’s intuition was confirmed by a series of attempts and experiments, that led to the conclusion that:

a Scheme for Classification, suited to books embodying multi-faceted subjects, should provide for two or more free facets and not merely one. This was the predisposing cause for designing a faceted Scheme for Classification (Ranganathan, 1965, pp. 13–14).

But there are exceptions to the rule! As well as in DDC the organization is not completely enumerative (Tartaglia, 1998, p. 7; 18-19), in CC the organization is not completely an analytic-synthetic one: for practical reasons – and for coherence and uniformity sake – a part of the three folded text of the CC, 6<sup>th</sup> edition (Ranganathan, 1963), is a scheme of classification. In fact, Part 3 contains Schedules of Classics and Sacred Books with special names, i.e. a list of ready-made class numbers for classic works in many main classes such as *L Medicine*, *Δ Spiritual Experience and mysticism*, *N Fine arts*, *O Literature*, *P Linguistics*, *Q Religion* and *R6 Indian Philosophy*.

These enumerative schedules of classics are an extremely useful tool, especially for Indian libraries. Nevertheless, it is very difficult to find some class numbers interesting for a Western library among the more than 126 pages of Part 3 of the

2 Translation from Italian by the author.

3 In *Elements of library classification*, the specific subject of a book is defined as “that division of knowledge whose extension and intension are equal to those of its thought-content” (Ranganathan, 1959, p. 4).

CC. If a Part 4 of CC devoted to Western classics – or at least literary classics – were available, the use of CC would be simpler, more efficient and more effective.

Nevertheless, the construction of this hypothetical Part 4 of CC would be a considerably time-consuming work, additionally to be done without any warrant that efforts would be rewarded by an actual and large use of new Schedules. Therefore, a question arises: is the automatic, quick, and cost-effective large-scale production of a Schedule of class numbers for Western literary classics possible? Moreover, the answer to this question implies considering a more general issue: which are the theoretical requirements for an automatic large-scale production of class numbers of a faceted scheme of classification?

Two prerequisites seem to be essential to answer in the affirmative both questions: 1) creation of class numbers must be a regular, exception-free and predictable process, so that it could be transformed into an algorithm; 2) modular, shareable, reusable and preferably free data needed to create facets provided for in the classification scheme must be available. The former prerequisite is granted by two specific characteristics of CC, i.e., the availability of the *Classic device* and the semantic approach provided by facet formula in Main Class *O Literature*. Both the characteristics allow defining an algorithm for creating new class numbers. The latter prerequisite is theoretically granted by Linked Open Data (LOD) available on the Semantic Web.

Are these prerequisites enough? The case study on *data.bnf.fr* and CC aims to develop the basis to evaluate the opportunity of a more extensive research on how LOD might be used for subject indexing by faceted classification schemes.

### **Classic device and semantics of Main Class *O Literature***

In CC, the classification process for literary works can be transformed in an algorithm thanks to classic device and facet formula of the Main Class *O Literature*. In CC, a classic work is ‘a book stimulating other books and literature on itself’; it is also defined as a ‘work (other than a Sacred Work or a Work of Literature) expounding some specialised subject, usually having embodiments in several versions, adaptations, and translations, inspiring other works on itself, and getting copied out and/or brought out in print even long after its origin’ (Ranganathan, 1963, p. 1.59, 1967, p. 486).<sup>4</sup> So, classic works exist in Literature, of course (e.g., Homer’s *Iliad*), but also in any other discipline: Darwin’s *The origin of the species*,

4 The definition is very clear and useful also in the IFLA LRM perspective, as it implies that a work should be considered a classic every time it stimulates either new expressions, or manifestations, or items, or other related works.

Aristotle's *Logic*, and Charles A. Cutter's *Rules* are all examples of classic works, as they stimulated other books and literature on themselves. Ranganathan notes:

This is strictly a classifier's definition. It naturally uses that quality of Classics which challenges a classificatory language to deal with it properly. It wants that a Classic and its associated literature should be arranged like a continuous spectrum, with nothing outside itself intervening. [...] The inherent qualities of a Classic stimulating such auxiliary literature are that: 1. It has elements of permanent value; 2. It is saturated with the personality of the author – which in itself was very powerful and highly organized; and 3. It is a seminal book cutting new ground, blazing new trail, stimulating new thought, and so on.

In CC two types of classic *devices* are provided; the former is devoted only to main classes *O Literature* and *Q Religion* and is based on a chronological principle, since chronological organization is considered more useful for these two disciplines. The latter applies to any other Main Class, and takes advantage of a subject approach, and subordinately of a bibliographical-literary approach. For example, the application of the classic device to the Aristotle's *Organon* results in an organization (*filiatory sequence*) of its expressions and manifestations, and subsequently of its comments (organized by their expressions and manifestations too).<sup>5</sup>

Main class *O Literature* is probably a *unicum* in CC, as its Schedules for classification consists in a list of less than 10 isolates (voices); it is also interesting because two out of four facets (Language and Time) of its formula can be obtained from common isolates Schedules.<sup>6</sup> Facet formula of Main Class *O Literature* is the following (CC, 2.94):

O [P], [P2] [P3], [P4]

where facet [P] represents the language, facet [P2] corresponds to literary form (see below), facet [P3] stands for author (represented by means of the birth year), and facet [P4] corresponds to the work facet (CC, 1.101-1.102). Every facet is intended to be used only when it applies.

Overall, facet formula means that, in Main Class *O Literature*, concepts and

5 For an example of application of the classic device to Aristotle's *Organon* and of the organization of expressions and manifestations of the work and its comments, see (Bianchini, 2012; Varghese, 2010).

6 Common Isolates in CC have their rough correspondent in Auxiliary Table of DDC.

punctuation must respect the following order:

O + Language + , + Literary Form + Author + , + Specific Work

The first facet [P] is determined by the Language Isolate schedule; Language facet expresses the language of the work (in a FRBR meaning).<sup>7</sup> For instance, as isolate numbers for English, German, and Russian languages are respectively 111, 113, and 142, English literature gets class number O111, German Literature gets O113, and Russian Literature gets O142. Facet formulas show that facet [P] is separated from facet [P2] by means of a comma. The second facet [P] is defined by the isolates of the Main Class *O Literature*; the array is fully presented in the *O Literature* schedule, and it is faceted by the literary form, as follows:

- 1 Poetry
- 2 Drama
- 3 Fiction, including short stories
- 4 Letters (literature written in the form of letters)
- 5 Oratorical
- 6 Other forms of prose
- 7 Champu<sup>8</sup>

Facet [P2] stands for the literary form of the work; its focus is “the form of literature which the book contains or about which the book treats” (CC, 1.98). For instance, a book containing Russian poetry takes the class number O142,1 (i.e. a comma and the isolate number 1 for Poetry are added to O142). In the same way, a book about English Drama gets class number O111,2.

Up to this point, CC and DDC facets are very similar: in both, the sequence of ideas and notation express first main class Literature, then language of the work, and lastly literary form of the work. Though, from the third facet onwards, CC becomes more specific, precise, and accurate, because it allows identification of both the author and, eventually, each single work of the author.

In the facet [P3], the author of the literary work or subject of the book is to be represented by the birth year. Years are determined by means of the Time isolate schedule, where each century is represented by a capital letter. For instance, the

7 In fact, in CC the other FRBR Group 1 entities – such as Expression, Manifestation, and Item – are represented by attributes expressed within the book number; see below.

8 Champu: a literary form of the Indian literary tradition (Kannada and Telugu languages) that uses a blend of prose and poetry.

letter J represents all the years from 1500 to 1599, N represents the range 1900-1999 and P represents the range 2000-2099. Accordingly, each single year of any century can be obtained by adding its two last figures to the Letter representing its century, so that 1564=J64, 1956=N56, and 2001=P01.

For example, by means of Schedule of Time, a book about Shakespeare's plays gets the class number O111,2J64 (as Shakespeare was born in 1564). The author indexing system created by Ranganathan is modern, because with this faceted approach a literary author is almost always given a unique notation.<sup>9</sup> It could also be used as an author identifier to some extent, and it allows representing an author both as a creator and as a subject in the same way within Ranganathan's classified catalogue.<sup>10</sup>

The fourth facet [P4] refers to the work or the works contained in the book, or about which the book is. The most relevant situation for the case study is that of books containing one work. In fact, in this case the isolate number for facet [P4] is a progressive special number assigned by the classifier,<sup>11</sup> following the chronological sequence of the author's works.<sup>12</sup>

To sum up, the reason why it is possible to produce CC class numbers by an algorithm is a specific characteristic of the main Class *O Literature*. Additionally, every data element to be used as a value in its facet formula – language, literary form, the author's birth year and the title of the work – is already available either in bibliographic and authority records or as Linked Open Data. For this reason, an algorithm for the automatic construction of class numbers for Main Class *O Literature* can check the values of the relevant data elements and transform them in the corresponding isolate numbers to synthesize the final class number of the book.

Data elements to be provided for the facet formula are the following:

- author's nationality, to identify the relevant national literature;
- literary form mostly associated with the author;
- author's birth year;
- number of works associated with the author;

9 It must be noted that the obtained class number represents more exactly *any* dramatic English author born in 1564, as O111,2 represents English Drama, and J64 stands just for 1564 (a year associated with William Shakespeare, but theoretically associable with any other dramatic author born in that year). Anyway, rules for disambiguation of authors are provided in case of ambiguity (CC, 1.100-101).

10 For further readings about Ranganathan's Classified Catalogue, see (Ranganathan, 1964).

11 The integer must be assigned by a special process to work out the exact amount and value of the figures (depending on the number of the author's works, if it is not exceeding 8 works, or it exceeds 8 but not 64 works, or it exceeds 64 works; CC, 1.101-102). For the sake of brevity, and provided that the process is completely automatable, a simply progressive integer was used in this study.

12 CC provides for the option to assign the isolate number for the work arbitrarily, if the chronological order is not easily ascertainable. In the case study, works can be ordered following either the tradition or general or specific reference works such as bibliographies, catalogues, thematic catalogues, and so on.

- date of each work associated with the author; and
- optionally, for a prospective Part 4 of the CC, the title of the works associated with the author.

Also, data elements correspondent to each facet must be available in a modular, shareable, and reusable form, to allow the large-scale production of class numbers for one or more national literatures.

## Literary data in the Semantic Web

The first attempt to use a freely available and structured source of linked data relating to literary works and authors was based on Wikipedia, and especially on its Wikidata Query Service.<sup>13</sup> Data are accessible by a SPARQL end point, to be operated by SPARQL Language but also by a simplified user interface. In Wikidata, data are recorded as entities and properties, and some properties are useful to extract bibliographic data relating to authors and works. Anyway, the simplest approach is to start from a list of author-entities, reconcile them with Wikidata and extract further data by means of their specific properties (for example, *given name* (P735), *date of birth* (P569), *occupation* (P106), *notable work* (P800), etc.).<sup>14</sup>

The sample of 200 Italian authors extracted by means of pages available in Wikipedia<sup>15</sup> highlighted a few issues preventing a successful reuse of Wikidata to create automatic class numbers in CC:

1. the percentage of literary authors in the data retrieved by property occupation (P106) as author (Q2500638) was relatively low (15%);
2. values for property occupation (P106) are not recorded consistently, as the string values ‘writer’, ‘author’, ‘poet’ and ‘novelist’ are not associated to the ‘creative person’ in a foreseeable and coherent way; nevertheless, they are essential to identify any literary creator;
3. property ‘author’ (P50) could not be used, as it can apply to any kind of work (such as a scientific paper, for example), and not just to literary works; additionally, property ‘creator’ (P170) was too generic.

13 <https://query.wikidata.org/>

14 The author is extremely grateful to Andrea Zanni for his helpfulness and aid in these steps and for the useful discussion about data available in Wikidata and their quality.

15 [https://it.wikisource.org/wiki/Categoria:Autori\\_del\\_XIX\\_secolo](https://it.wikisource.org/wiki/Categoria:Autori_del_XIX_secolo).

In general, the analysis of the sample data extracted from Wikidata showed interesting and useful data when relating to a single author, but also data lacking consistency, regularity and completeness when considered with respect to the universe of literary authors. Consequently, it was impossible to extract data relevant to all the authors and works of a specific national literature, or period, etc. The failed approach to extract coherent and complete data from Wikidata suggested finding a different, specialized source of linked open data; the prospective source had to contain authoritative data, created from the scratch and compliant with cataloguing standards. After an exploratory and quick analysis on both *data.bnf.fr* and *datos.bne.es*,<sup>16</sup> the data available from Bibliothèque nationale de France were chosen. Data necessary to produce class numbers for a sample of authors were extracted from *data.bnf.fr* by its SPARQL end point (<http://data.bnf.fr/sparql/>) (see Query in Appendix 1).

## The sample

As the main goal of the case study was to evaluate the feasibility of the process, a restricted data sample was identified, collected and elaborated. Italian authors were chosen, according to the following requirements:

- authors born between 1800 and 1900;
- ‘Italy’ as country associated to the person;<sup>17</sup>
- ‘Italian’ as language associated to the person;<sup>18</sup>
- ‘Literature’ (literal) as field of activity of the person.<sup>19</sup>

For each matched author, the following data were extracted from the SPARQL end point (Figure 1):

- URL of the work (column 1)
- Title of the work (column 2)
- Language of the work (column 3)
- URI of the date of the work (column 4)

16 *datos.bne.es* showed a major issue: date of works was not available, while it was mandatory for the chronological arrangement of the works and their subsequent numbering. For this reason, *data.bnf.fr* was chosen as the source for the data sample.

17 Value <http://id.loc.gov/vocabulary/countries/it> for property “rdagroup2elements:countryAssociatedWithThePerson”.

18 Value <http://id.loc.gov/vocabulary/iso639-2/ita> for “rdagroup2elements:languageOfThePerson”.

19 Value ‘Littératures’ for property “rdagroup2elements:fieldOfActivityOfThePerson”.

- Date of the work (in numeric format) (column 5)
- Author URI (column 6)
- Full name of the author (in the form “Name Surname”; column 7)
- Field “Note” (usually containing a non-structured description of the author; (column 8)
- Given name of the author (column 9)
- Family name of the author (column 10)
- Birth date of the author (column 11)
- Death date of the author (column 12)

id	titre	date	type	nom	nom complet	note	premiere	debut	fin	debut	fin
http://data.bnf.fr/ark:/12148/cb115668227about	"Pai de Tolosan"@it	1816	Salvatore Commensano	Salvatore Commensano	Peintre, poète, dramaturge et librettiste"	Salvatore"	Commensano"	1801	1812		
http://data.bnf.fr/ark:/12148/cb146882727about	"L'avis Miller"@it	1849	Salvatore Commensano	Salvatore Commensano	Peintre, poète, dramaturge et librettiste"	Salvatore"	Commensano"	1801	1812		
http://data.bnf.fr/ark:/12148/cb141669060about	"Don Pasquale"@it	1847	Giovanni Raffain"	Giovanni Raffain"	"Romancier et librettiste"	Giovanni"	Raffain"	1807	1811		
http://data.bnf.fr/ark:/12148/cb141669060about	"Don Pasquale"@it	1847	Giovanni Raffain"	Giovanni Raffain"	"Écrivain suédois en anglais"	Giovanni"	Raffain"	1807	1811		
http://data.bnf.fr/ark:/12148/cb12707618about	"I ministri de Firenze"@it	1817	Carlo Colloidi"	Carlo Colloidi"	"A nous traduit du français en italien"	Carlo"	Colloidi"	1826	1890		
http://data.bnf.fr/ark:/12148/cb12707618about	"I ministri de Firenze"@it	1817	Carlo Colloidi"	Carlo Colloidi"	"Écrivain et journaliste - Fondateur du journal de satire politique "Il lampione" (1849-1853), repris sous le titre "La scaramuccia" à partir de 1860 - Auteur des "Aventures de Panocchio", parues pour la première fois en 1880 dans le "Giornale dei bambini" - Traducteur de Charles Perrault"	Carlo"	Colloidi"	1826	1890		
http://data.bnf.fr/ark:/12148/cb12707618about	"Le avventure di Panocchio"@it	1883	Carlo Colloidi"	Carlo Colloidi"	"A nous traduit du français en italien"	Carlo"	Colloidi"	1826	1890		
http://data.bnf.fr/ark:/12148/cb12707618about	"Le avventure di Panocchio"@it	1883	Carlo Colloidi"	Carlo Colloidi"	"Écrivain et journaliste - Fondateur du journal de satire politique "Il lampione" (1849-1853), repris sous le titre "La scaramuccia" à partir de 1860 - Auteur des "Aventures de Panocchio", parues pour la première fois en 1880 dans le "Giornale dei bambini" - Traducteur de Charles Perrault"	Carlo"	Colloidi"	1826	1890		
http://data.bnf.fr/ark:/12148/cb1462992727about	"Frattelli d'Italia"@it	1847	Goffredo Mameli"	Goffredo Mameli"	"Poète et patriote, l'un des fondateurs de la république romaine, auteur d'Hymnes patriotiques dont l'un des textes l'Hymne national italien en 1848"	Goffredo"	Mameli"	1827	1849		
http://data.bnf.fr/ark:/12148/cb1427092174about	"Le condizioni di un italiano"@it	1867	Ippolito Nievo"	Ippolito Nievo"	"Romancier"	Ippolito"	Nievo"	1811	1861		

FIGURE 1 HTML output of the query on the *data.bnf.fr* SPARQL endpoint

Automatic collection of data was difficult with reference to one data element: data element ‘field of activity’, which was queried for ‘Literature’, was not available for an unquantifiable but relevant part of Italian authors recorded in *data.bnf.fr*. For this reason, the sample of Italian authors collected by the query resulted not to be complete. The data sample was of 31 authors and 97 literary works of the Italian Literature of the XIX and XX centuries. The total amount of authors and works included in the data sample are very far from the completeness; nevertheless, the sample was enough for the purpose of the case study.

## Sample data processing

Sample data were exported in CSV format and imported in OpenRefine, a ‘powerful tool for working with messy data: cleaning it; transforming it from one format into another; and extending it with web services and external data.’<sup>20</sup> Data

20 <http://openrefine.org/>

processing required the following steps:

- 1) deduplication of records on the basis of identity of both Author's URI and URI of work;
- 2) arrangement by the author, date of the work and title of the work, in order to obtain the sequence required by CC work arrangement;
- 3) for each author, progressive numbering of his/her works arranged as per point 2;
- 4) clustering of authors for literary genre by data element Description (i.e. field 'Note' in *data.bnf.fr*).

Step no. 4 was the most difficult and longest one and required a small manual data processing by clustering function in OpenRefine. In fact, a relevant point to obtain correct CC class numbers in the main class *O Literature* was to identify each author's literary genre, as all the works by an author are classified by means of the author class number, which depends on the literary form associated with the author most. Information relating to the author's literary form were contained in the field "Note", but in an unstructured form (Figure 1; Column 8 "Tipo\_Autore"). Nevertheless, the field 'Note' was useful even if unstructured, because it allowed to identify the prevailing literary form with the first characteristic used to describe the author: for instance, if the author's description was 'Romancier et auteur dramatique', Literary form facet in the CC class number was assumed to be "Romancier" as the right and relevant value and the author was associated to Novel as preferred literary form. To collocate author forms manually but more quickly, the value of the textual string containing the non-structured description was truncated to 20 characters. Just three authors (9%) required a specific manual classification, as their main literary role was unclear after the truncation.<sup>21</sup>

After cleaning data relating to literary form for each author, data were restructured automatically through the following steps:

columns were rearranged and selected to respect the order and the value required by facet formula;

verbalstring values were translated in their equivalent notation, by means of string transformation utilities available in OpenRefine, so that three new col-

21 In particular: Carlo Collodi was associated with the truncated string 'A aussi traduit du f', and Giovanni Ruffini was associated to the string 'Ecrivait aussi en an', both phrases without any reference to a literary form; Giorgio De Chirico was associated to the string 'Peintre, dessinateur' (a clear string but not usable to the literary classification). During the manual data revision Collodi and Ruffini were associated to Novel form, while De Chirico was excluded from the sample of Italian authors (even if he is associated to 843.91 in OCLC Classify; <http://classify.oclc.org/classify2/>).

umns (FaccettaForma, FaccettaAutore, and NumeroOpera) were automatically (see Figure 2, columns nos. 3, 5 and 7);

by rearrangement of authors' names (arranged alphabetically by family names), work title (arranged chronologically and, in suborder, alphabetically), and the sum of the values of the three newly created columns, a class number for each literary work was created in OpenRefine;

the CC class number tables for the complete sample of 31 Italian literary authors and 97 literary works were produced by the extraction of Full name of the author, + Work Title + Class number from OpenRefine to a text document (Appendix 1).

Facet / Filter	Undo / Redo 77	32 records		Extensions: Freebase							
Extract... Apply...		Show as: rows records	Show: 5 10 25 50 records	« first < previous 1 10 next > last »							
Filter:		NomeCompleto	tipo_autore	FaccettaForma	nato	FaccettaAutore	titolo_opera	NumeroOpera	NumeroCC	autore	opera
70. Remove 1 rows		2. Camillo Boito	Critico	6	1836	M36	Senso, nuove storielle vane	1	O121,3M36,1	http://data.bnf.fr/ark:/12148/cb119026346#about	http://data.bnf.fr/ark:/12148/cb12207618c#about
71. Create new column FaccettaForma based on column tipo_autore by filling 99 rows with grel:value.replace(/Poeta/, "1");replace(/Drammaturgo/, "2");replace(/Narratore/, "3");replace(/Critico/, "6")		3. Camillo Scabarbo	Poeta	1	1888	M88	Trucchi	1	O121,1M88,1	http://data.bnf.fr/ark:/12148/cb115454527#about	http://data.bnf.fr/ark:/12148/cb115454527#about
72. Create new column FaccettaAutore based on column nato by filling 99 rows with grel:value.replace(/9/, "M")		4. Carlo Colodi	Narratore	3	1826	M26	I matieri de Firenze	1	O121,3M26,1	http://data.bnf.fr/ark:/12148/cb118973660#about	http://data.bnf.fr/ark:/12148/cb12207618c#about
73. Move column titolo_opera to position 6		5. Carlo Dossi	Narratore	3	1849	M49	Le avventure di Pinocchio	2	O121,3M26,2	http://data.bnf.fr/ark:/12148/cb118973660#about	http://data.bnf.fr/ark:/12148/cb123129155#about
74. Move column titolo_opera to position 5		6. Carlo Emilio Gadda	Narratore	3	1893	M93	La colonia felice	1	O121,3M49,1	http://data.bnf.fr/ark:/12148/cb120202017#about	http://data.bnf.fr/ark:/12148/cb11900025c#about
75. Move column NumeroOpera to position 7					1849	M49	Note azzurre	2	O121,3M49,2	http://data.bnf.fr/ark:/12148/cb120202017#about	http://data.bnf.fr/ark:/12148/cb11900025c#about
76. Move column NumeroOpera to position 6					1893	M93	L'Adalgisa	1	O121,3M93,1	http://data.bnf.fr/ark:/12148/cb119020165#about	http://data.bnf.fr/ark:/12148/cb11900025c#about
77. Create new column NumeroCC based on column NumeroOpera by filling 99 rows with grel:"O121,3" + "+" + colref["FaccettaForma"] value + colref["FaccettaAutore"] value + "+" + colref["NumeroOpera"] value					1893	M93	Quer pasticciaccio brutto de via Merulana	2	O121,3M93,2	http://data.bnf.fr/ark:/12148/cb119020165#about	http://data.bnf.fr/ark:/12148/cb11900025c#about
					1893	M93	La cognizione del dolore	3	O121,3M93,3	http://data.bnf.fr/ark:/12148/cb119020165#about	http://data.bnf.fr/ark:/12148/cb11900025c#about
					1893	M93	I Luigi di Francia	4	O121,3M93,4	http://data.bnf.fr/ark:/12148/cb119020165#about	http://data.bnf.fr/ark:/12148/cb11900025c#about

FIGURE 2 Data cleaning and transformation in OpenRefine

## Discussion

The case study and its results put forward some interesting points with respect to the twofold starting issue, that is both whether an automatic, quick, affordable and large-scale creation of classification tables of CC class number for the literature is possible and, from a more general approach, which are the requirements for developing an automatic process for the creation of faceted class numbers. In the case study, the creation of class numbers for isolates of the main class *O Literature* relating to single works of specific authors resulted workable. Nevertheless, two possible issues emerged from the experimentation:

- during the first steps of the process, the extraction of linked open data on literary authors resulted partial, because only a part of the potentially relevant authors was associated with the tag 'literature'; and,

- during the step of identification of the relevant literary class for each author, data processing was not completely automated because the source data from *data.bnf.fr* were not structured (for example by means of a controlled vocabulary), but they were available only in a non-structured form.

Both issues can be bypassed in a relatively simple way, that is, widening the harvest of LOD to more than one namespace and reconciling data from many namespaces. These means should allow a more exact identification of literary authors and getting different attributes for the authors from various name spaces, for instance getting the author's name, surname, birth and death dates, and works from one source and author's preferred literary form and classification from another source. A specific feature of OpenRefine permits saving the whole data transformation process in a JSON file, that can be reused and applied to newly extracted data (for instance, a second sample of Italian literature authors, or the full extraction of French literature authors, and so on). The availability of a routine process assures replicability of the process to a larger scale, especially if controlled and high-quality data could be downloaded from many LOD sources.

As to the more general issue, the first thing to be observed is a prerequisite for the classification scheme. It must be based on regular semantics and syntaxes – that is to say on facets and facet formula, to permit class numbers to be obtained by an algorithm. This prerequisite prevents in a relevant way the applicability of the described process to classification schemes widely spread all over the word but not structured in such a way (e.g. Dewey Decimal Classification). An example of this is that the case study purposely took in exam the specific facet formula provided for single works of literary authors in the CC. In fact, Special Table in Part 3 of CC deals with these cases, and not with other kind of publications such as selections or collections of one or more authors. Moreover, the success of the process requires availability of semantic data with specific characteristics. It must be noted that:

- semantic data availability alone is not sufficient, as in the case of Wikidata, where property 'occupation' (P106) is provided for but data values are recorded inconsistently (i.e. lacking consistency and standardization provided for ICP (IFLA Section on Cataloguing Code, 2016)); so that, even if Wikidata is potentially usable as LOD, the way its data are produced and recorded did not assure sufficient consistence to their reuse in the case study;
- semantic data availability in one repository could not be sufficient alone, even if data are produced and recorded consistently; in order for data to

be reusable for the described process, their completeness with respect to the values of any facet in the facet formula is necessary. For example, at *datos.bne.es*, data relating to which and how many works were linked to an author were not accessible; and,

- complete and coherent semantic data availability is not sufficient alone, because their completeness must be checked also in relation to their domain (for example, Literature); indeed, many Italian literature authors were not found by the query at *data.bnf.fr* SPARQL endpoint because of the lacking of value ‘Littératures’ in the property “field of activity of the person”; so that, accessible data allowed to perform the process, but results are likely to be quantitatively poor.

Therefore, for a more structured project aiming at large scale class number production a carefully designed data harvest is mandatory in relation to the number and the quality and completeness of data elements.

### Five lines of future research

The case study suggests a few hypotheses for further research. First, the issue of the difficult determination of the prevailing literary form associated with an author based on the unstructured biographical description obtained from *data.bnf.fr* could be fixed by the individuation of the third figure of DDC class number mainly associated with that author; for instance, an author mainly associated to DDC class number 853.xx could be associated to Italian novel. At least, till more reliable sources with this kind of data will be available.

Another possible line of research should include the individuation of other facet formulas within the CC scheme, and their analysis in order to determine the feasibility of such a process of automatic production of class numbers in other disciplinary fields. An example of such facet formulas in the CC are those relating to the classic works in other main classes, to be described by means of their relationships to other works (derived works, works based on works, and so on). Another example is that of main class *NR Music*, which is based on facets like space [P1], time [P2], a list of form of music publications [P3],<sup>22</sup> and instruments [M].<sup>23</sup>

<sup>22</sup> For instance, Word/Libretto, Notation, Form, Keeping time, Dramatic music, Orchestral music.

<sup>23</sup> For instance, *Wind instrument* (that includes *Pipe, Flute, Organ* etc.), *Stringed instrument* (that includes *Vina, Violin, Piano*, etc.).

A third interesting field of study could be the application of LOD to the production of book numbers, that is the part of classified call numbers that is used to distinguish among, and to identify, resources belonging in the same class. Ranganathan designed a facet formula for book numbers that is particularly interesting for many reasons: 1) it allows a unique identification of each book of a classified collection; 2) it is based on data elements usually available in MARC family records; 3) its application produces a sequence of resources that is fully compliant with a FRBR-zed arrangement of Expressions, Manifestations, and Items of a Work, by using core characteristics of each FRBR entity.<sup>24</sup> Since the identification of each item of a set of resources sharing the same class number is an issue not limited to CC but to any available scheme of classification, this line of research could be useful in a wide range of situations.

A fourth possible study could be directed to investigate and detect the presence of facet formulas in other classification schemes (for example, in class 400 Linguistics or in class 800 Literature of the Dewey Decimal Classification, or Universal Decimal Classification), to search further possible applications of the process.

The last, most complex and ambitious research goal is the design of a purely faceted classification scheme based mainly, if not exclusively, on data usually and easily available as linked open data. In many classification schemes, tables for common concepts such as language, space and time are provided. These concepts are largely and widely represented, recorded, and available as linked open data in the semantic web. They are a fundamental and relevant starting point for such a project.

## REFERENCES

- BIANCHINI, C. (2012). Arrangement of FRBR Entities in Colon Classification Call Numbers. *Cataloging & Classification Quarterly*, 50(5-7), 1-21, <https://doi.org/10.1080/01639374.2012.679877>.
- BIANCHINI, C. (2017). *Book number: uno strumento per l'organizzazione delle collezioni. Manuale ad uso dei bibliotecari*. Milano: Editrice Bibliografica.
- IFLA Section on Cataloguing Code. (2016). *Statement of International Catalogu-*

<sup>24</sup> Core characteristics for the book number facet formula are *language* (referred to as the Expression) and *date of publication* (relating to the Manifestation). For further readings on the book number and its relationships with FRBR, see (Bianchini, 2012), and, more recently (in Italian), (Bianchini, 2017, pp. 73-94).

- ing Principles (ICP)*. A. Galeffi, M.V. Bertolini, R.L. Bothmann, E. Escolano Rodríguez and D. McGarry (eds.). Den Haag: IFLA.
- RANGANATHAN, S.R. (1959). *Elements of library classification : Based on lectures delivered at the University of Bombay ... 1944, and in the schools of librarianship in Great Britain ... 1956*. B. I. Palmer (ed.) (2nd ed.). London: The Association of Assistant Librarians.
- RANGANATHAN, S.R. (1963). *Colon classification. Basic classification* (6th ed. rev). Bangalore: Sarada Ranganathan Endowment for Library Science.
- RANGANATHAN, S.R. (1964). *Classified catalogue code with additional rules for Dictionary catalogue code* (5th ed.). Bombay: Asia Publishing House.
- RANGANATHAN, S.R. (1965). *A descriptive account of the Colon Classification*. Bangalore: Sarada Ranganathan Endowment for Library Science.
- RANGANATHAN, S.R. (1967). *Prolegomena to Library Classification* (3rd ed.). Bombay: Asia Publishing House.
- TARTAGLIA, S. (1998). *Ordine di citazione e principio di faccettazione nella Classificazione Decimale Dewey*. Udine: Forum.
- VARGHESE, M. (2010). Relevance of a Classified Catalogue in the FRBR perspective and a proposed Model with ISBD descriptions and Faceted Class number as Key attribute. *Cataloging & Classification Quarterly*, 46(3), 281-304.

## APPENDIX NO.1

### SPECIAL TABLES. ITALIAN LITERATURE

This Appendix contains CC class numbers of Italian literature works obtained processing linked open data extracted from *data.bnf.fr*.

Titles of the works just indicate the work; they neither correspond to the titles on the first edition nor represent a 'uniform title'; in fact, identification is to be obtained by CC work number. As all the works are classics, form of their titles varies depending on the reference source (for instance, *Enciclopedia Treccani*, BNCf opac, Wikipedia). Therefore, the form of the title is taken from *Dizionario letterario Bompiani delle opere e dei personaggi di tutti i tempi e di tutte le letterature* (Milano, Bompiani, 1964-1972), when available. Otherwise, it is taken directly from *data.bnf.fr*.

Author	Work	Work class number
Corrado Alvaro	<i>Uomo è forte</i>	O121,3M95,1
	<i>Letà breve</i>	O121,3M95,2
	<i>Memorie del mondo sommerso</i>	O121,3M95,3
	<i>Lunga notte di Medea</i>	O121,3M95,4
Luigi Antonelli	<i>La donna in vetrina</i>	O121,2M77,1
Camillo Boito	<i>Senso</i>	O121,6M36,1
Anna Bonacci	<i>L'ora della fantasia</i>	O121,2M92,1
Salvatore Cammarano	<i>Pia de' Tolomei</i>	O121,2M01,1
	<i>Luisa Miller</i>	O121,2M01,2
Dino Campana	<i>Canti orfici</i>	O121,1M85,1
Carlo Collodi	<i>I misteri di Firenze</i>	O121,3M26,1
	<i>Le avventure di Pinocchio</i>	O121,3M26,2
Edmondo De Amicis	<i>La vita militare : bozzetti</i>	O121,1M46,1
	<i>Cuore</i>	O121,1M46,2
	<i>Primo maggio</i>	O121,1M46,3
Grazia Deledda	<i>Canne al vento</i>	O121,3M71,1
	<i>Il segreto dell'uomo solitario</i>	O121,3M71,2
Emilio De Marchi	<i>Demetrio Pianelli</i>	O121,3M51,1
Federico De Roberto	<i>La sorte</i>	O121,3M61,1
	<i>I viceré</i>	O121,3M61,2
Carlo Dossi	<i>La colonia felice</i>	O121,3M49,1
	<i>Note azzurre</i>	O121,3M49,2
Carlo Emilio Gadda	<i>L'Adalgisa</i>	O121,3M93,1
	<i>Quer pasticciaccio brutto de via Merulana</i>	O121,3M93,2
	<i>La cognizione del dolore</i>	O121,3M93,3
	<i>I Luigi di Francia</i>	O121,3M93,4
	<i>La meccanica</i>	O121,3M93,5
	<i>Meditazione milanese</i>	O121,3M93,6
Corrado Govoni	<i>Le fiale</i>	O121,1M84,1
Guido Gozzano	<i>La via del rifugio</i>	O121,1M83,1
	<i>I colloqui</i>	O121,1M83,2

Author	Work	Work class number
Luigi Illica	<i>Andrea Chénier</i>	O121,2M57,1
Goffredo Mameli	<i>Fratelli d'Italia</i>	O121,1M27,1
Filippo Tommaso Marinetti	<i>Mafarka, il futurista</i>	O121,1M76,1
	<i>Parole in libertà</i>	O121,1M76,2
	<i>Gli indomabili</i>	O121,1M76,3
	<i>Teatro della sorpresa</i>	O121,1M76,4
	<i>Cinque sintesi radiofoniche</i>	O121,1M76,5
Eugenio Montale	<i>Ossi di seppia</i>	O121,1M96,1
	<i>Il balcone dentro alle occasioni</i>	O121,1M96,2
	<i>Le occasioni</i>	O121,1M96,3
	<i>Mediterraneo</i>	O121,1M96,4
	<i>Mottetti</i>	O121,1M96,5
	<i>Finisterre</i>	O121,1M96,6
	<i>Farfalla di Dinard</i>	O121,1M96,7
	<i>Languilla</i>	O121,1M96,8
	<i>La bufera e altro</i>	O121,1M96,9
	<i>Satura</i>	O121,1M96,10
	<i>Diario postumo</i>	O121,1M96,11
Dario Niccodemi	<i>Prete Pero</i>	O121,2M74,1
Ippolito Nievo	<i>Confessioni di un Italiano</i>	O121,3M31,1
Luigi Pirandello	<i>L'esclusa</i>	O121,2M67,1
	<i>Il turno</i>	O121,2M67,2
	<i>Il fu Mattia Pascal</i>	O121,2M67,3
	<i>L'umorismo</i>	O121,2M67,4
	<i>I vecchi e i giovani</i>	O121,2M67,5
	<i>Quaderni di Serafino Gubbio operatore</i>	O121,2M67,6
	<i>La giara</i>	O121,2M67,7
	<i>Il giuoco delle parti</i>	O121,2M67,8
	<i>Sei personaggi in cerca d'autore</i>	O121,2M67,9
	<i>Enrico IV</i>	O121,2M67,10
	<i>Novelle per un anno</i>	O121,2M67,11
	<i>Ciascuno a suo modo</i>	O121,2M67,12
	<i>Uno, nessuno e centomila</i>	O121,2M67,13

Author	Work	Work class number
	<i>Lazzaro</i>	O121,2M67,14
	<i>Questa sera si recita a soggetto</i>	O121,2M67,15
	<i>I giganti della montagna</i>	O121,2M67,16
Giovanni Ruffini	<i>Don Pasquale</i>	O121,3M07,1
Umberto Saba	<i>Canzoniere</i>	O121,1M83,1
	<i>Autobiografia</i>	O121,1M83,2
	<i>Scorciatoie e raccontini</i>	O121,1M83,3
Camillo Sbarbaro	<i>Trucioli</i>	O121,1M88,1
Renato Serra	<i>Esame di coscienza di un letterato</i>	O121,6M84,1
Scipio Slataper	<i>Il mio carso</i>	O121,3M88,1
Italo Svevo	<i>Una vita</i>	O121,3M61,1
	<i>Senilità</i>	O121,3M61,2
	<i>La coscienza di Zeno</i>	O121,3M61,3
	<i>Una burla riuscita</i>	O121,3M61,4
	<i>La novella del buon vecchio e della bella fanciulla</i>	O121,3M61,5
Giuseppe Tomasi di Lampedusa	<i>Il gattopardo</i>	O121,3M96,1
	<i>La sirena</i> <sup>25</sup>	O121,3M96,2
Giuseppe Ungaretti	<i>Il porto sepolto</i>	O121,1M88,1
	<i>Sentimento del tempo</i>	O121,1M88,2
	<i>Allegria</i>	O121,1M88,3
Giovanni Verga	<i>I Carbonari della montagna</i>	O121,3M40,1
	<i>Sulle lagune</i>	O121,3M40,2
	<i>Storia di una capinera</i>	O121,3M40,3
	<i>Rosso Malpelo</i>	O121,3M40,4
	<i>Vita dei campi</i>	O121,3M40,5
	<i>I Malavoglia</i>	O121,3M40,6
	<i>Novelle rusticane</i>	O121,3M40,7
	<i>Cavalleria rusticana</i>	O121,3M40,8
	<i>Mastro don Gesualdo</i>	O121,3M40,9
	<i>Don Candeloro e C.i</i>	O121,3M40,10
	<i>La caccia alla volpe</i>	O121,3M40,11
	<i>Dal tuo al mio</i>	O121,3M40,12

# OD SEMANTIČKOG WEBA DO FACETNE KLASIFIKACIJE STUDIJA SLUČAJA I PET PRAVACA BUDUĆEG ISTRAŽIVANJA

## KLJUČNE RIJEČI:

*facetna klasifikacija, semantički web, Klasifikacija s dvotočkom, data.bnf.fr, klasifikacija talijanske književnosti*

## SAŽETAK

Opisana je studija slučaja automatskog stvaranja brojeva facetne klasifikacije za posebne teme, koje polazi od podataka dostupnih na semantičkom webu. Cilj je istražiti dva pitanja: 1) Je li moguće automatski, brzo, ekonomično i u velikom opsegu generirati brojeve Klasifikacije s dvotočkom za književne klasičke? i 2) Općenito: Koji su uvjeti za automatsko generiranje brojeva facetne klasifikacijske sheme? Kakva su istraživanja potrebna za razvoj takvog procesa? Studija slučaja usredotočena je na podatke izvučene iz data.bnf.fr i na njihovu ponovnu uporabu u cilju dobivanja točnih i potpunih brojeva Klasifikacije s dvotočkom za uzorak talijanskih književnika (rođenih od 1800. do 1900.). Predstavljen je postupak kreiranja klasifikacijskih brojeva koji polazi od identifikacije i ekstrakcije podataka te njihovog čišćenja, transformacije i prijevoda u klasifikacijsku notaciju pomoću alata Open Refine. Raspravlja se o rezultatima studije slučaja i njihovoj važnosti za početna pitanja te se identificiraju istraživačka pitanja za moguće buduće pravce razvoja.