

ANÁLISIS DE UN DATASET

Para la práctica de JSON y JSON Schema se pide inspeccionar un dataset de Kaggle o de Github y escribir un mini-reportaje analizándolo.

De entre las opciones existentes, se ha escogido el dataset “books.json” del repositorio GitHub <https://github.com/ozlerhakan/mongodb-json-files.git> . El dataset en cuestión consta de 431 libros del mundo de la informática y la programación.

Una imagen en crudo de este fichero .json se puede ver a continuación:

```

1 [{"id": 1, "title": "Android in Action, Second Edition", "isbn": "1935182722", "pageCount": 416, "publishedDate": "2011-01-14T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
2 [{"id": 2, "title": "Specification by Example", "isbn": "1971200884", "pageCount": 80, "publishedDate": "2011-06-03T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
3 [{"id": 3, "title": "Flex 3 in Action", "isbn": "1933987648", "pageCount": 576, "publishedDate": "2009-02-07T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
4 [{"id": 4, "title": "Flex 4 in Action", "isbn": "1933988312", "pageCount": 425, "publishedDate": "2008-10-10T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
5 [{"id": 5, "title": "Collective Intelligence in Action", "isbn": "1933988312", "pageCount": 425, "publishedDate": "2008-10-10T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
6 [{"id": 6, "title": "Zend Framework in Action", "isbn": "1933988320", "pageCount": 432, "publishedDate": "2008-12-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
7 [{"id": 7, "title": "Flex on Java", "isbn": "1933988797", "pageCount": 265, "publishedDate": "2010-15-08T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
8 [{"id": 8, "title": "Griffon in Action", "isbn": "1935182234", "pageCount": 375, "publishedDate": "2010-06-04T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
9 [{"id": 9, "title": "Flexible Java", "isbn": "1933988599", "pageCount": 592, "publishedDate": "2009-01-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
10 [{"id": 10, "title": "Flexible Rails", "isbn": "1933988599", "pageCount": 592, "publishedDate": "2009-01-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
11 [{"id": 11, "title": "Hello! Flex 4", "isbn": "1933988762", "pageCount": 258, "publishedDate": "2009-11-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
12 [{"id": 12, "title": "Coffeehouse", "isbn": "1884777384", "pageCount": 316, "publishedDate": "1997-07-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
13 [{"id": 13, "title": "Team Foundation Server 2010 in Action", "isbn": "1935188592", "pageCount": 344, "publishedDate": "2010-12-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
14 [{"id": 14, "title": "PongBd in Action", "isbn": "1935182870", "pageCount": 80, "publishedDate": "2011-12-17T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
15 [{"id": 15, "title": "Distributed Application Development with PowerBuilder 6.0", "isbn": "1884777686", "pageCount": 584, "publishedDate": "1998-06-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
16 [{"id": 16, "title": "Jaguar Developer with PowerBuilder 7", "isbn": "1884777864", "pageCount": 559, "publishedDate": "1999-08-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
17 [{"id": 17, "title": "Timing Jaguar", "isbn": "1884777686", "pageCount": 362, "publishedDate": "2000-07-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
18 [{"id": 18, "title": "Hibernate in Action", "isbn": "193294153X", "pageCount": 400, "publishedDate": "2006-08-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
19 [{"id": 19, "title": "Hibernate in Action (Chinese Edition)", "isbn": "193294153X", "pageCount": 400, "publishedDate": "2006-08-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
20 [{"id": 20, "title": "Java Persistence with Hibernate", "isbn": "1932941885", "pageCount": 880, "publishedDate": "2006-11-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
21 [{"id": 21, "title": "JSTL in Action", "isbn": "1930105299", "pageCount": 480, "publishedDate": "2002-07-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
22 [{"id": 22, "title": "Designing Hard Software", "isbn": "133040493X", "pageCount": 350, "publishedDate": "1997-02-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
23 [{"id": 23, "title": "Designing Hard Software", "isbn": "133040493X", "pageCount": 350, "publishedDate": "1997-02-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
24 [{"id": 24, "title": "Hibernate Search in Action", "isbn": "1933988649", "pageCount": 488, "publishedDate": "2008-12-21T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
25 [{"id": 25, "title": "Query in Action", "isbn": "1933988355", "pageCount": 376, "publishedDate": "2008-01-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
26 [{"id": 26, "title": "Query in Action, Second Edition", "isbn": "1935182323", "pageCount": 488, "publishedDate": "2010-06-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
27 [{"id": 27, "title": "Ruby for Rails", "isbn": "1932940499", "pageCount": 532, "publishedDate": "2006-05-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
28 [{"id": 28, "title": "The Well-Grounded Rubyist", "isbn": "1933988657", "pageCount": 520, "publishedDate": "2009-04-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
29 [{"id": 29, "title": "ASP.NET 4.0 in Practice", "isbn": "1935182463", "pageCount": 584, "publishedDate": "2011-05-15T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
30 [{"id": 30, "title": "PFC Programmer's Reference Manual", "isbn": "1884777554", "pageCount": 368, "publishedDate": "1998-06-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
31 [{"id": 31, "title": "Graphics File Formats", "isbn": "133040454", "pageCount": 484, "publishedDate": "1995-06-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
32 [{"id": 32, "title": "Visual Object Oriented Programming", "isbn": "133040454", "pageCount": 280, "publishedDate": "1995-02-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
33 [{"id": 33, "title": "IOS in Practice", "isbn": "1017291269", "pageCount": 325, "publishedDate": "2011-11-01T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
34 [{"id": 34, "title": "Silverlight 2 in Action", "isbn": "1933988428", "pageCount": 480, "publishedDate": "2008-10-31T00:00:00-0700", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
35 [{"id": 35, "title": "The Quick Python Book, Second Edition", "isbn": "193518220X", "pageCount": 360, "publishedDate": "2010-01-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
36 [{"id": 36, "title": "Internet and Intranet Applications with PowerBuilder 6", "isbn": "1884777680", "pageCount": 330, "publishedDate": "2010-12-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"},
37 [{"id": 37, "title": "Practical Methods for Your Year 2000 Problem", "isbn": "188477752X", "pageCount": 236, "publishedDate": "1998-01-01T00:00:00-0800", "thumbnailUrl": "https://s3.amazonaws.com/AKIA5CSRIADUWPRP"}]

```

Para poder realizar el análisis es necesario tener estos datos en un formato más legible y, por ello, se ha escogido un libro en concreto y se han adaptado sus propiedades:

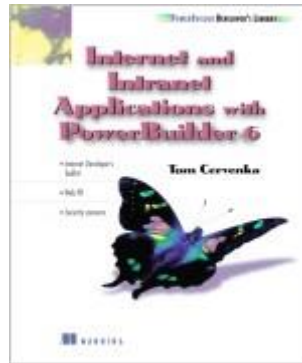
```

1  {
2    "id" : 46,
3    "title" : "Internet and Intranet Applications with PowerBuilder 6",
4    "isbn" : "1884777600",
5    "pageCount" : 390,
6    "publishedDate" : { "$date" : "2000-12-01T00:00:00.000-0800" },
7    "thumbnailUrl" : "https://s3.amazonaws.com/AKIAJC5RLADLWVRPFDQ.book-thumb-images/cervenka.jpg",
8    "longDescription" : "If you're a PowerBuilder programmer, Internet and Intranet Applications with PowerBuilder 6 is your ticket to learning Web.
9    PB and related technologies. The book covers everything you need to know to build web browser and server programs with the PowerBuilder 6 Internet Toolkit.
10   Also covered is how to write winsock programs with PB, and Distributed PB is covered to the extent necessary to learn Web.PB.",
11    "status" : "PUBLISH",
12    "authors" : [ "Tom Cervenka" ],
13    "categories" : [ "PowerBuilder" ]
14  }

```

Como se puede observar, las propiedades de cada libro cubren las características típicas, desde el título y los autores hasta el ISBN. Además, cuenta con un identificador unívoco para cada libro "id".

Los tipos de cada una de ellas son variados: hay strings, numbers, arrays y colecciones. Un ejemplo de colección sería la propiedad “publishDate”, la cuál es definida mediante una clave-valor como es la variable “\$date” y la fecha en formato string. Por otro lado, la propiedad “thumbnailUrl” tiene como tipo un string que abstrae la URL de la imagen de la portada del libro.



Por último, cabe destacar los tipos arrays en propiedades como “authors” o “categories”, pues al poder tener varios valores es adecuado definirlas así. Concretamente en este ejemplo el libro solo tiene una categoría y un autor, pero en varios de los otros libros del dataset sí que tienen más.

Analizando el resto de los libros del conjunto se puede ver que en muchos de ellos está, como en el analizado en este reporte, la propiedad “longDescription”; en otros, en vez de esta, aparece la propiedad “shortDescription”; y, en otros, no aparece ninguna de las dos mencionadas. Esto lleva a concluir que en el esquema de este JSON esa propiedad en concreto no aparecerá como requerida mientras que el resto si lo estarán.