Práctica 8 - Información Semántica

Hugo Pérez Fernández. U0250708@uniovi.es

Sistemas de Información para la Web. Grado de Ingeniería Informática. EII. universidad de Oviedo. Campus de los Catalanes. Oviedo.

1 Introducción

A continuación se realizarán los ejercicios de la practica 8 de la asignatura de Sistemas de Información para la Web, estos se dividirán en dos conjuntos, en el primero se trabajara sobre una noticia (Anexo.4.1) para marcarla con microdata y se resolveran unas cuestiones sobre dicha tecnología, y en el segundo se realizaran una serie de modificaciones en un JSONLD dado para mejorar la información que provea.

2 Asignación de entidades y creación de microdatos.

Se trata de seleccionar un extracto de una noticia y realizar la asignación de entidades a este en forma de microdatos mediante las herramientas Dandelion, Schema.org y Prueba de datos Estructurados de Google.

2.1 Definición de entidades mediante Dandelion.

Se ha seleccionado el extracto (Anexo.4.1) de la noticia sobre el ransomware que ha afectado a empresas de España (Noticia aquí). Y, usando la herramienta Dandelion, se han obtenido las siguientes entidades:

- Organisations:
 - Cadena SER.
 - Prisa Radio.
 - Everis.
 - Radio Madrid.
- Concepts:
 - Ransomware.
 - Virus Informático.
 - Internet.
 - Computadora personal.
 - Prensa escrita.
 - Medio de Comunicación.
 - Empresa.

2.2 Refinamiento de entidades con Schema.org

Una vez que tenemos un primer aproximamiento a las entidades que aparecen en el extracto (Anexo.4.1) mediante el uso de la herramienta Dandelion, ahora pasaremos a completar la lista generada mediante nuestro conocimiento del dominio de la noticia y el apoyo de Schema.org, de manera que localicemos mas items o entidades y sus propiedades.

Tras realizar esto generamos el siguiente html con microdatos:

```
<div itemscope itemtype="http://schema.org/NewsArticle">
1
2
       <h1 itemprop="headline">Everis y Prisa Radio sufren
           un grave ciberataque que
3
            secuestra sus sistemas</hl>
4
       <span itemprop="datePublished" content="2019-11-04">
           2019-11-04 < / span >
       <span itemprop="articleBody">
5
            Varias empresas españolas han sufrido hoy un
               serio ataque de
           <span itemscope itemtype="http://schema.org/</pre>
 7
               Software Application">
8
           <span itemprop="name">Ransomware
9
           </span>que recuerda al vivido a mediados de 2017
10
           <span itemscope itemtype="http://schema.org/</pre>
               Software Application">
11
           <span itemprop="name">Wannacry</span>
12
           </span>. Los primeros ataques confirmados de
               forma oficial
13
            los han sufrido
14
           <span itemscope itemtype="http://schema.org/</pre>
               Organization">
15
           <span itemprop="name">La Cadena SER</span>
16
           </span> y otras emisoras de <span itemscope
               itemtype="http://schema.org/Organization">
17
           <span itemprop="name">Prisa Radio
18
           </span>, pero también varias consultoras
            tecnológicas, de las cuales < span itemscope
19
               itemtype="http://schema.org/Organization">
20
           <span itemprop="name">Everis</span>
21
           </span> ha confirmado oficialmente estar afectada
               . Estamos sufriendo
22
           un ataque masivo de virus a la red de <span
               itemscope itemtype="http://schema.org/
               Organization">
23
           <span itemprop="name">Everis</span>
24
           </span>. Por favor, mantengan los <span itemscope
                itemtype="http://schema.org/Product">
```

```
25
           <span itemprop="name">PC</span>
26
           <img itemprop="image" src="">
27
           <span itemprop="description"></span>
28
           </span> apagados. Es el
29
           mensaje interno que ha remitido <span itemscope
               itemtype="http://schema.org/Organization">
           <span itemprop="name">Everis</span> a sus
30
               empleados, según ha podido
31
            confirmar este diario
32
           y han publicado también varios medios
               especializados. La compañía
33
            confirma que ha enviado a sus
            trabajadores a casa hasta que puedan solventar la
34
                incidencia.
35
36
           <span itemscope itemtype="http://schema.org/</pre>
               Organization">
37
           <span itemprop="name">Prisa Radio/span> ha
               sufrido esta
38
           madrugada un ataque de <span itemscope itemtype="
               http://schema.org/SoftwareApplication">
           <span itemprop="name">Virus</span>
39
40
           </span> que ha tenido una afectación grave y
            generalizada de todos nuestros <span itemscope
41
               itemtype="http://schema.org/Thing">
42
           <span itemprop="name">sistemas informáticos</span</pre>
43
           </span>. Los técnicos especializados en este tipo
44
           de situaciones aconsejan encarecidamente la
               desconexión total de
45
            todos los sistemas con el fin
46
            de evitar la propagación del <span itemscope
               itemtype="http://schema.org/
               Software Application">
47
           <span itemprop="name">Virus</span>. Hablamos, por
                tanto, de
48
           una situación de extrema emergencia,
49
           ha comunicado esta mañana la empresa en un
               mensaje interno
50
            al que ha tenido acceso este diario.
51
52
           A partir de ese momento, se irá chequeando puesto
                a puesto
53
           -siguiendo las instrucciones precisas
```

```
54
           que os enviará el departamento de Sistemas-para
               autorizar
55
           en los casos en que haya garantía
            absoluta la puesta en marcha de cada equipo. Por
56
               el momento,
           y hasta nuevo aviso, la emisión de
57
           <span itemscope itemtype="http://schema.org/</pre>
58
               Organization">
59
           <span itemprop="name">La Cadena SER</span>
60
           </span> queda centralizada en <span itemscope
               itemtype="http://schema.org/Organization">
61
           <span itemprop="name">Radio Madrid
62
           </span>; quedan anuladas todas las emisiones
               locales y
63
            regionales, salvo aquellas que hayan sido
               autorizadas
64
            expresamente por la <span itemscope itemtype="
               http://schema.org/Organization">
65
           <span itemprop="name">Dirección de Antena/span>
66
           </span>.
           En este momento, la seguridad es la máxima de la
67
               compañía
68
           por encima de cualquier otro compromiso.
69
            Cualquier acción individual no autorizada puede
               poner en
70
            peligro el trabajo de rescate que se está
71
            llevando a cabo, cierra el mensaje.
72
           </span>
73
   </div>
```

Listing 1.1. Archivo HTML con Microdatos embebidos.

Al validar el documento html anterior mediante la herramienta para pruebas de datos estructurados de Google genera el siguiente resultado:

Como se ve en la Fig.1 al validar el html generado este presenta errores y advertencias en las diferentes entidades, esto se debe a decisiones tomadas ante los incovenientes que presentan los microdatos, que se explicarán a continuación.

La entidad NewsArticle presenta los siguientes problemas:

- Errores:

- Falta el atributo Author obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo image obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo publisher obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.
- Advertencias:

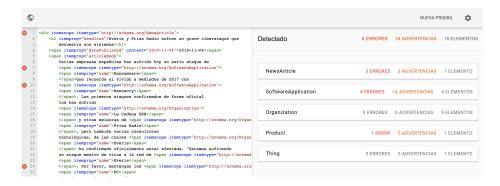


Fig. 1. Resultado del testeo del html con microdatos embebidos.

- Falta el atributo DateModified recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo mainEntityOfPage recomendado: Este es debido a que el extracto de la noticia no contiene esa información.

Las 4 entidades SoftwareAplication presentan los mismos problemas:

- Errores:

• En este caso el error es que son necesarios al menos dos de los atributos recomendados que se indican como advertencias a continuación.

Advertencias:

- Falta el atributo aggregateRating recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo applicationCategory recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo offers recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo operatingSystem recomendado: Este es debido a que el extracto de la noticia no contiene esa información.

La entidad Product presenta los siguientes problemas:

– Errores:

- Falta el atributo offers obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo review obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo aggregateRating obligaotrio: Este es debido a que el extracto de la noticia no contiene esa información.

- Advertencias:

• Falta el atributo brand recomendado: Este es debido a que el extracto de la noticia no contiene esa información.

- Falta el atributo description recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo image recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta el atributo sku recomendado: Este es debido a que el extracto de la noticia no contiene esa información.
- Falta algun atributo de tipo identificador global recomendado: Este es debido a que el extracto de la noticia no contiene esa información.

2.3 Cuestiones sobre el uso de microdatos.

Ya realizado el proceso de los apartados anteriores queda plantearse las siguientes preguntas:

¿Deberían tener valor para itemid todos los items? ¿Qué valor debería asignarse?

Depende, todos los itemid deberian tener un itemid de manera que se pudiera identificar cada uno de ellos unequivocamente, pero no en todos los casos, puesto que sobrecargaría los textos. Además esto favorecería la propiedad de la web Semántica por la que toda entidad debe ser enlazable.

¿Si hubiera varios candidatos cuál escogerías? ¿Por qué?.

En el caso de que hubiera varios candidatos para el identificador de un item elegiria un valor alfanumerico, como realiza WikiData por ejemplo, de este modo evitariamos colisiones y problemas de idiomas.

¿Cuál crees que es la solución de compromiso que podría resultar menos polémica en la mayoría de casos?.

Para decidir si hay que añadir un itemid o no la solucion seria si ese item tiene un identificador, ya que quizas no lo tenga, y si mejora la información que nos da.

¿ Qué inconvenientes concretos te ha supuesto la obligación de incrustarlos metadatos allí donde te "forzaba" la estructura del texto?.

El mayor problema que ha supesto incrustar los metadatos en la estructura del texto es que, por el modo en el que se hace cualquiera de ellos debe aparecer en el texto de manera que algunas propiedades obligatorias de las entidades quedan sin rellenar porque no aparecen en el texto, y esto al revisarlo mediante la herramienta de Google da errores, que si quisieramos solucionar la unica manera posible para ello seria modificar el texto que se fuera a mostrar para que apareciesen todos los datos de la propiedades de cada entidad.

3 JSON-LD

A continuación se presentará la ultima version del jsonld que se facilita en la practica (Anexo.4.2), el cual se ha modificado de la siguiente forma:

 Se han añadido a las entidades Periodical para que el id ahora sea la web oficial de la publicaciones.

- Se ha añadido al contexto del jsonld la url de DOI "https://doi.org/", y se han modificado las los items de tipo Chapter para que el id apunte a un "digital object identifier.
- Se ha asociado al id de la entidad Person de Daniel Gayo su web oficial. Y de igual modo con las entidades Universidad de Oviedo y Departamento de informatica con sus id's y las web's oficiales.
- Se han añadido a los items de tipo Periodical el atributo hasPart para señalar el articulo del que se habla de toda la publicacion.

1	<div></div>
2	
3	Varias empresas españolas han sufrido hoy un serio ataque de 'ransomware' que recuerda al
4	vivido a mediados de 2017 con Wannacry. Los primeros ataques confirmados de forma oficial
5	los han sufrido la Cadena SER y otras emisoras de Prisa Radio, pero también varias consultoras
6	tecnológicas, de las cuales Everis ha confirmado oficialmente estar afectada. Estamos sufriendo
7	un ataque masivo de virus a la red de Everis. Por favor, mantengan los PCs apagados. Es el
8	mensaje interno que ha remitido Everis a sus empleados, según ha podido confirmar este diario
9	y han publicado también varios medios especializados. La compañía confirma que ha enviado a sus
10	trabajadores a casa hasta que puedan solventar la incidencia.
11	
12	
13	PRISA Radio ha sufrido esta madrugada un ataque
	de virus que ha tenido una afectación grave y
14	generalizada de todos nuestros sistemas informáticos. Los técnicos especializados en este tipo
15	de situaciones aconsejan encarecidamente la desconexión total de todos los sistemas con el fin
16	de evitar la propagación del virus. Hablamos, por tanto, de una situación de extrema emergencia
17	ha comunicado esta mañana la empresa en un mensaje interno al que ha tenido acceso este diario.
18	

```
19
       >
           A partir de ese momento, se irá chequeando puesto
20
                a puesto -siguiendo las instrucciones
               precisas
21
            que os enviará el departamento de Sistemas- para
               autorizar en los casos en que haya garantía
22
            absoluta la puesta en marcha de cada equipo. Por
               el momento, y hasta nuevo aviso, la emisión de
           Cadena SER queda centralizada en Radio Madrid;
23
               quedan anuladas todas las emisiones locales y
24
            regionales, salvo aquellas que hayan sido
               autorizadas expresamente por la Dirección de
               Antena.
25
           En este momento, la seguridad es la máxima de la
               compañía por encima de cualquier otro
               compromiso.
26
            Cualquier acción individual no autorizada puede
               poner en peligro el trabajo de rescate que se
27
            llevando a cabo, cierra el mensaje.
28
       29 < / \operatorname{div} >
  <script type="application/ld+json">
31
   {
32
       "@context": {
         "cc": "http://creativecommons.org/ns#",
33
34
         "ctag": "http://commontag.org/ns#",
         "dc": "http://purl.org/dc/elements/1.1/",
35
36
         "dc1": "http://purl.org/dc/terms/",
         "dc11": "http://purl.org/dc/elements/1.1/",
37
         "dcat": "http://www.w3.org/ns/dcat#"
38
39
         "dcterms": "http://purl.org/dc/terms/",
40
         "doi": "https://doi.org/",
         "foaf": "http://xmlns.com/foaf/0.1/".
41
42
         "gr": "http://purl.org/goodrelations/v1#",
         "grddl": "http://www.w3.org/2003/g/data-view#",
43
         "hcalendar": "http://microformats.org/profile/
44
             hcalendar#",
         "hcard": "http://microformats.org/profile/hcard#",
45
         "ical": "http://www.w3.org/2002/12/cal/icaltzd#",
46
         "ma": "http://www.w3.org/ns/ma-ont#",
47
         "md": "http://www.w3.org/ns/md#",
48
         "og": "http://ogp.me/ns#"
49
50
         "org": "http://www.w3.org/ns/org#",
51
         "owl": "http://www.w3.org/2002/07/owl#",
```

```
52
          "prov": "http://www.w3.org/ns/prov#",
53
         "qb": "http://purl.org/linked-data/cube#",
          "rdf": "http://www.w3.org/1999/02/22-rdf-syntax-ns#
54
         "rdfa": "http://www.w3.org/ns/rdfa#",
55
         "rdfs": "http://www.w3.org/2000/01/rdf-schema#",
56
          "rev": "http://purl.org/stuff/rev#",
57
          "rif": "http://www.w3.org/2007/rif#",
58
         "rr": "http://www.w3.org/ns/r2rml#",
59
         "schema": "http://schema.org/",
60
         "sd": "http://www.w3.org/ns/sparql-service-
61
             description#",
62
         "sioc": "http://rdfs.org/sioc/ns#",
63
         "skos": "http://www.w3.org/2004/02/skos/core#",
64
          "skosxl": "http://www.w3.org/2008/05/skos-xl#",
65
         "v": "http://rdf.data-vocabulary.org/#",
         "vcard": "http://www.w3.org/2006/vcard/ns#",
66
         "void": "http://rdfs.org/ns/void#",
67
         "wdr": "http://www.w3.org/2007/05/powder#",
68
         "wdrs": "http://www.w3.org/2007/05/powder-s#"
69
70
         "wdsr": "http://www.w3.org/2007/05/powder-s#",
          "xhv": "http://www.w3.org/1999/xhtml/vocab#",
71
72
          "xsd": "http://www.w3.org/2001/XMLSchema#"
73
       "@graph": [
74
75
76
            "@id": "_: N87d20637af2445ec83dea37c29d9085d",
           "@type": "schema: Chapter",
77
78
            "schema:isPartOf": {
              "@id": "doi:10.1017/CBO9781316182635"
79
80
            },
81
            "schema:name": "Political Opinion"
82
83
            "@id": "https://ieeexplore.ieee.org/xpl/
84
               RecentIssue.jsp?punumber=93",
            "@type": "schema: Periodical",
85
            "schema:name": "IEEE Multimedia",
86
87
           "hasPart": "doi:10.1109/MMUL.2015.47"
88
89
            "@id": "https://www.emerald.com/insight/
90
               publication/issn/1066-2243",
91
            "@type": "schema: Periodical",
92
            "schema:name": "Internet Research",
```

```
93
             "hasPart": "issn:1066-2243"
94
           },
95
             "@id": "https://www.computer.org/csdl/magazine/ic
96
97
             "@type": "schema: Periodical",
             "schema:name": "IEEE Internet Computing",
98
             "hasPart": "doi:10.1109/MIC.2012.137"
99
100
101
             "@id": "https://cacm.acm.org",
102
             "@type": "schema: Periodical",
103
             "schema: name": "Communications of the ACM",
104
             "hasPart":"doi:10.1145/2001269.2001297 |"
105
106
107
             "@id": "_{-}: Nb6734e117ddf48a4a3f95b95e5239145",
108
             "@type": "schema: Organization",
109
             "schema:name": "Department of Computer Science"
110
111
112
             "@id": "http://danigayo.info",
113
             "@type": "schema:Person",
114
             "schema: affiliation": {
115
               "@id": "_:Nd46c7cbe934541d68dc746e686f6354c"
116
117
118
             "schema:jobTitle": "associate professor",
             "schema:name": "Daniel Gayo-Avello"
119
120
121
             "@id": "urn:isbn:9781107500075",
122
             "@type": "schema:Book",
123
             "schema:name": "\"Twitter: A Digital Socioscope\"
124
125
             "schema: publisher": "Cambridge University Press"
126
127
             "@id": "http://www.uniovi.es",
128
             "@type": "schema: CollegeOrUniversity",
129
130
             "schema:department": {
               "@id": "http://www.di.uniovi.es"
131
132
             "schema:name": "University of Oviedo"
133
134
           }
135
```

```
136 }
137 </script>
```

Listing 1.2. Archivo HTML con JSONLD embebido.

4 Anexos

4.1 Extracto de la noticia seleccionada.

Varias empresas españolas han sufrido hoy un serio ataque de 'ransomware' que recuerda al vivido a mediados de 2017 con Wannacry. Los primeros ataques confirmados de forma oficial los han sufrido la Cadena SER y otras emisoras de Prisa Radio, pero también varias consultoras tecnológicas, de las cuales Everis ha confirmado oficialmente estar afectada. "Estamos sufriendo un ataque masivo de virus a la red de Everis. Por favor, mantengan los PCs apagados". Es el mensaje interno que ha remitido Everis a sus empleados, según ha podido confirmar este diario y han publicado también varios medios especializados. La compañía confirma que ha enviado a sus trabajadores a casa hasta que puedan solventar la incidencia.

"PRISA Radio ha sufrido esta madrugada un ataque de virus que ha tenido una afectación grave y generalizada de todos nuestros sistemas informáticos. Los técnicos especializados en este tipo de situaciones aconsejan encarecidamente la desconexión total de todos los sistemas con el fin de evitar la propagación del virus. Hablamos, por tanto, de una situación de extrema emergencia", ha comunicado esta mañana la empresa en un mensaje interno al que ha tenido acceso este diario.

"A partir de ese momento, se irá chequeando puesto a puesto -siguiendo las instrucciones precisas que os enviará el departamento de Sistemas- para autorizar en los casos en que haya garantía absoluta la puesta en marcha de cada equipo. Por el momento, y hasta nuevo aviso, la emisión de Cadena SER queda centralizada en Radio Madrid; quedan anuladas todas las emisiones locales y regionales, salvo aquellas que hayan sido autorizadas expresamente por la Dirección de Antena. En este momento, la seguridad es la máxima de la compañía por encima de cualquier otro compromiso. Cualquier acción individual no autorizada puede poner en peligro el trabajo de rescate que se está llevando a cabo", cierra el mensaje.

4.2 JSONLD original

```
1  {
2     "@context": {
3         "cc": "http://creativecommons.org/ns#",
4         "ctag": "http://commontag.org/ns#",
5         "dc": "http://purl.org/dc/elements/1.1/",
6         "dc1": "http://purl.org/dc/terms/",
```

```
7
       "dc11": "http://purl.org/dc/elements/1.1/",
8
       "dcat": "http://www.w3.org/ns/dcat#",
9
       "dcterms": "http://purl.org/dc/terms/",
       "foaf": "http://xmlns.com/foaf/0.1/",
10
       "gr": "http://purl.org/goodrelations/v1#",
11
        "grddl": "http://www.w3.org/2003/g/data-view#",
12
        "hcalendar": "http://microformats.org/profile/
13
           hcalendar#",
       "hcard": "http://microformats.org/profile/hcard#",
14
15
       "ical": "http://www.w3.org/2002/12/cal/icaltzd#",
16
       "ma": "http://www.w3.org/ns/ma-ont#",
17
        "md": "http://www.w3.org/ns/md#",
       "og": "http://ogp.me/ns#".
18
        "org": "http://www.w3.org/ns/org#",
19
       "owl": "http://www.w3.org/2002/07/owl#",
20
21
       "prov": "http://www.w3.org/ns/prov#",
       "qb": "http://purl.org/linked-data/cube#",
22
       "rdf": "http://www.w3.org/1999/02/22-rdf-syntax-ns#", "rdfa": "http://www.w3.org/ns/rdfa#",
23
24
       "rdfs": "http://www.w3.org/2000/01/rdf-schema#",
25
26
       "rev": "http://purl.org/stuff/rev#",
27
       "rif": "http://www.w3.org/2007/rif#",
        "rr": "http://www.w3.org/ns/r2rml#",
28
29
       "schema": "http://schema.org/",
30
        "sd": "http://www.w3.org/ns/sparql-service-
           description#",
31
       "sioc": "http://rdfs.org/sioc/ns#",
32
       "skos": "http://www.w3.org/2004/02/skos/core#",
33
       "skosxl": "http://www.w3.org/2008/05/skos-xl#",
       "v": "http://rdf.data-vocabulary.org/#",
34
       "vcard": "http://www.w3.org/2006/vcard/ns#",
35
36
       "void": "http://rdfs.org/ns/void#",
37
       "wdr": "http://www.w3.org/2007/05/powder#",
        "wdrs": "http://www.w3.org/2007/05/powder-s#"
38
39
       "wdsr": "http://www.w3.org/2007/05/powder-s#",
       "xhv": "http://www.w3.org/1999/xhtml/vocab#",
40
        "xsd": "http://www.w3.org/2001/XMLSchema#"
41
     },
"@graph": [
42
43
44
          "@id": "_: N87d20637af2445ec83dea37c29d9085d",
45
          "@type": "schema: Chapter",
46
          "schema:isPartOf": {
47
            "@id": "urn:isbn:9781107500075"
48
49
          },
```

```
50
          "schema:name": "Political Opinion"
51
52
          "@id": "_: N237308aadd894e9cb88f3f9bd3ca7f85",
53
         "@type": "schema: Periodical",
54
          "schema:name": "IEEE Multimedia"
55
56
57
         "@id": "_: N0f982a4b306847a78b9496ab615d9502",
58
         "@type": "schema: Periodical",
59
          "schema:name": "Internet Research"
60
61
62
         "@id": "_{-}: N57a8d3cded6d4061af9d3cb1e7acd39f",
63
         "@type": "schema:Periodical",
64
65
          "schema:name": "IEEE Internet Computing"
66
67
          "@id": "\_: N8381ebd04bda429fbb2c1d900d8b0e2c"\;,
68
69
          "@type": "schema: Periodical",
70
          "schema:name": "Communications of the ACM"
71
72
73
          "@id": "_:Nb6734e117ddf48a4a3f95b95e5239145",
         "@type": "http://schema.org/Organization",
74
          "http://schema.org/name": "Department of Computer
75
             Science"
76
77
         "@id": " \_: N8013fd042aca458b9f5373abc2c9bbda" \, ,
78
          "@type": "schema:Person",
79
80
          "schema: affiliation": {
            "@id": "_:Nd46c7cbe934541d68dc746e686f6354c"
81
82
83
          "schema: jobTitle": "associate professor",
          "schema:name": "Daniel Gayo-Avello"
84
85
86
          "@id": "urn:isbn:9781107500075",
87
         "@type": "schema:Book",
88
          "schema:name": "\"Twitter: A Digital Socioscope\"",
89
         "schema: publisher": "Cambridge University Press"
90
91
92
          "@id": "_:Nd46c7cbe934541d68dc746e686f6354c",
93
```

Listing 1.3. JSON Original facilitado en la practica.