

Bases de Datos II

Francisco Javier Mercader Martínez

Índice

1	Recuperación de datos y formatos de serialización	1
1.1	Necesidad de formatos de serialización	1
1.2	Características a analizar	1
1.2.1	Pandas	2
1.2.2	XML (eXtensible Markup Language)	2
1.2.3	CSV (Comma-Separated Values)	6
1.2.4	JSON (JavaScript Object Notation)	8
1.2.5	Apache Avro	10

Tema 1: Recuperación de datos y formatos de serialización

1.1) Necesidad de formatos de serialización

- Los formatos de **serialización** son vitales para el intercambio de datos
- Más en el ámbito *Big Data* (atacan a la "V" de la *variabilidad*)
- A lo largo de los años se han diseñado formatos de serialización
- Algunos son **más eficientes** que otros.
- Algunos se han **estandarizado**.
- En cualquier caso, son fundamentales para transmitir información, ya sea a través de **ficheros de disco** o bien para **comunicación por red**.
- La mayoría de los formatos son **de propósito general**, por lo que una misma información se puede codificar **de varias formas** (por ejemplo, como INSERT en SQL, ficheros CSV, etc.)
- Nuestro objetivo es conocer los formatos más usados para elegir el correcto en cada ocasión.

1.2) Características a analizar

- **¿Es un protocolo estándar?**- Con ello nos referimos a si está avalado por cuerpo de estándares o su especificación
- **¿Permite la codificación binaria?**- A la hora de transmitir grandes cantidades de información (ya sea en forma almacenada o bien a través de la red) es fundamental un protocolo binario para ahorrar espacio/tiempo.
- **¿Es legible por los humanos?**- A veces es interesante poder depurar un protocolo por parte de un humano. Esta idea comenzó con protocolos como XML, pero se ha ido abandonando porque no resulta muy factible salvo en ocasiones muy específicas.
- **¿Soporta referencias?**- A veces tenemos que relacionar partes de un conjunto de datos. Es interesante que los mecanismos de serialización permitan referenciar otras partes de un documento o de una comunicación.
- **¿Su estructura está definida por un Esquema o IDL?**- Al igual como sucedía con el lenguaje DDL de SQL, a veces es interesante que los datos sean conformes a algún esquema, también llamado IDL (*Interface Definition Language*).
- **¿Es extensible?**- A veces es necesario acomodar datos que no se ajustan estrictamente a un esquema, o que son directamente no-estructurados.
- **¿Poseen un API estandarizado?**- Si los formatos de serialización poseen un API estandarizado será más sencillo no sólo compartir los datos, sino también compartir el código de serialización/deserialización (también llamado **marshalling**)

- [Tabla resumen](#)

	Estándar	Bin?	Humano?	Ref?	IDL?	Ext?	API?
Apache Avro	Sí	Sí	No	N/A	Sí (acoplado)	Sí	N/A
CSV	Parcial (RFC4180)	No	Sí	No	No	Parcial	No
JSON	Sí (RFC7159)	No (BSON)	Sí	Sí (RFC6901)	Parcial	Sí	No
Thrift	No	Sí	Parcial	No	Sí (acoplado)	Sí	No
XML	Sí	Parcial	Sí	Sí	Sí	No	No

1.2.1) Pandas

- Pandas es una librería open source construida sobre Numpy.
- Permite una preparación, limpieza y análisis rápido de los datos.
- Una de sus principales características es la de visualización de datos.
- Puede trabajar con una gran variedad de fuentes musicales.
- Para instalar pandas es tan sencillo como ejecutar una de dichas instrucciones en el terminal

```
pip install pandas
conda install pandas
```

[Dataframe](#)

- La herramienta más conocida y usada de Pandas son los DataFrames.
- Permite almacenar datos tabulares en dos dimensiones similar a una hoja de cálculo o una base de datos relacional.
- Las columnas de datos

```
1 df = pd.DataFrame(randn(5,4),index='A B C D E'.split(),columns='W X Y Z'.split())
```

	W	X	Y	Z
A	-1.040684	-1.692150	1.707399	-1.257771
B	-0.403809	-1.024655	2.060558	-0.242150
C	-0.856354	0.173779	1.124053	-0.434952
D	0.282316	-1.349518	-0.076797	1.077644
E	-0.152517	-0.603708	-0.812906	0.807102

1.2.2) XML (eXtensible Markup Language)

- Meta-lenguaje de etiquetas derivado de SGML.
- Motivación:
 - Intercambio de datos en Internet
- Reúne los requisitos de un lenguaje de intercambio de información:
 - Simple: al estar basado en etiquetas y legible
 - Independiente de la plataforma: codificación UNICODE
 - Estándar y amplia difusión: W3C
 - Definición de estructuras complejas: DTD, Schemas
 - Validación y transformación: DTD, XSLT
 - Integración con otros sistemas
- Facilita procesamiento lado cliente.
- **No muy utilizado para BigData. Ha perdido tracción, es demasiado complejo finalmente y es muy poco eficiente en cual al porcentaje datos/metadatos**

[Ejemplo](#)

```

1 <?xml version="1.0" encoding="ISO-8859-1" ?>
2 <DISCO CODIGO="B000067FSG">
3   <TITULO>Estrella de Mar</TITULO>
4   <ARTISTA>Amaral</ARTISTA>
5   <ESTILO>Pop</ESTILO>
6   <REFERENCIA>
7     <EDITORIA>Virgin</EDITORIA>
8     <AÑO_EDICION>2002<AÑO_EDICION>
9   </REFERENCIA>
10  <MUSICOS>
11    <MUSICO ROL="cantante">Amaral</MUSICO>
12    <MUSICO ROL="guitarra">Juan Aguirre</MUSICO>
13  </MUSICOS>
14 </DISCO>

```

- Instrucciones de procesamiento (línea 1)
- Raíz (línea 2)
- Etiquetas y atributos

• DTD

```

1 <!ELEMENT DISCO (TITULO, ARTISTA, ESTILO?, REFERENCIA, MUSICOS)>
2 <!ATTLIST DISCO CODIGO ID #REQUIRED>
3 <!ATTLIST DISCO TIPO=(CD | LP | DVD) "CD">
4 <!ELEMENT TITULO (#PCDATA)>
5 ...
6 <!ELEMENT REFERENCIA (EDITORIA, AÑO_EDICION) >
7 <!ELEMENT MUSICOS (MUSICO*)>
8 ...

```

Describe los documentos XML \Rightarrow **Validación**

• DTD (ii)

Documento XML:

- **Válido:** sigue la estructura de un DTD

```

1 <!DOCTYPE web-app
2   PUBLIC "-//Sun Microsystems, Inc.//DTD Web Application 2.2//EN"
3   "http://java.sun.com/j2ee/dtds/web-app_2_2.dtd">

```

- **Bien formado:** Sigue las reglas de XML

Limitaciones:

- No es XML
- Tipado limitado de datos
- No soporta espacios de nombres

- Familias de estándares

- Schemas

- Mismo propósito que un DTD, pero con mayor riqueza semántica.
- Sintaxis basada en XML.
- Contiene tipos predefinidos
- Espacios de nombres (*namespaces*)

```
1 <?xml version="1.0">
2 <xsd:schema xmlns:xsd="http://www.w3c.org/2000/08/XMLSchema">
3 <xsd:element name="Disco" type="DiscoTipo"/>
4 <xsd:complexType name="DiscoTipo">
5   <xsd:attribute name="codigo" type="String"/>
6   <xsd:sequence>
7     <xsd:attribute name="Titulo" type="String"/>
8     <xsd:attribute name="Artista" type="String"/>
9     <xsd:attribute name="Referencia" type="ReferenciaTipo"/>
10    ...
```

- NameSpaces:

- Espacios de nombres para cualificar elementos y atributos evitando la colisión de nombres.
- `xmlns:xsd="http://www.w3c.org/2000/08/XMLSchema"`

- XSLT:

- Definición de reglas de transformación de documentos

- XSL:

- Definición de hojas de estilos.

- XPath:

- Para hacer referencia a partes de un documento
- `/DISCO[Titulo="Estrella de Mar"]`, `/DISCO//MUSICOS[1]`

- XLink:

- Enlace documentos entre sí.

- XPointer:

- Enlace de secciones dentro de un documento.

- XQuery:

- Consultas XML.

- **Parsers**

- API SAX:

- Acceso secuencial al documento

- Modelo de programación basado en eventos (*callbacks*)
- Simple y rápido: consume pocos recursos.
- Sólo consulta.
- API DOM:
 - Construye una estructura arbórea a partir del documento
 - Potente, pero más costoso.
 - Permite actualizaciones.
 - Ideal para estructuras complejas.
- En Python, existen numerosas formas de poder leer documentos XML.

• API SAX

Librería **xml.sax**

```

1 import xml.sax
2
3 class XMLHandler(xml.sax.ContentHandler):
4     def __init__(self):
5         # Inicializamos variables de interés
6
7         # Se llama cuando comienza un nuevo elemento
8         def StartElement(self, tag, attributes):
9             pass
10        # Se llama cuando un elemento acaba
11        def endElement(self, tag):
12            pass
13
14 parser = xml.sax.make_parser()
15 parser.setFeature(xml.sax.handler.feature_namespaces, 0)
16 Handler = XMLHandler()
17 parser.setContentHandler(Handler)
18 parser.parse('models.xml') # nombre del documento a analizar

```

Cómo podríamos procesar con **xml.sax** este documento

```

1 <collection shelf="New Arrivals">
2     <model number="ST001">
3         <price>35000</price>
4         <qty>12</qty>
5         <company>Samsung</company>
6     </model>
7     <model number="RW345">
8         <price>46500</price>
9         <qty>14</qty>
10        <company>Onida</company>
11    </model>

```

```
12 </collection>
```

- Parser DOM

Librería `xml.dom`

```
1 # Procesar un determinado fichero
2 file = minidom.parse('model.xml')
3
4 # Obtener los elementos con un determinado tag
5 modelos = fil.getElementsByTagName('modelo')
6
7 # Obtener el atributo 'nombre' del segundo
  modelo
8 print('modelo #2 atributos:')
9 print(modelos[1].attributes['nombre'].value)
10
11 # El datos de un item específico
12 print('\nmodelo #2 datos:')
13 print(modelos[1].firstChild.data)
```

```
1 <data>
2     <modelos>
3         <modelo name='modelo1'>
4             modelo1abc
5         </modelo>
6         <modelo name="modelo2">
7             modelo2abc
8         </modelo>
9     </modelos>
10 </data>
```

- Librería BeautifulSoup

```
1 from bs4 import BeautifulSoup
2
3 # Leemos el fichero
4 with open('models.xml', 'r') as f:
5     data = f.read()
6
7 # Pasamos los datos al parse
8 bs_data = BeautifulSoup(data, 'xml')
9
10 # Buscamos todas las instancias 'unique'
11 b_unique = bs_data.find_all('unique')
12 print(b_unique)
13
14 # Usamos .find búsquedas más concretas
15 b_name = bs_data.find('child', {'name': 'Acer'})
16 print(b_name)
```

```
1 <modelo>
2     <child name="Acer" qty="12">
3         Portátil Acer
4     </child>
5     <unique>
6         Número de modelo
7     </unique>
8     <child name="Acer" qty="7">
9         Exclusive
10    </child>
11    <unique>
12        1.200€
13    </unique>
14 </modelo>
```

1.2.3) CSV (Comma-Separated Values)

- Formato basado en columnas normalmente separado por coma
- Sin embargo, el formato admite variaciones:
 - Con o sin cabecera con los nombres de las columnas
 - Separador de columnas (comas, tabuladores, etc.)
 - Codificación de caracteres (UTF-8, latin-1, etc.)

- Escapado de caracteres (por ejemplo, una comilla doble como dos comillas dobles(" ") ó como \")
- Comillas opcionales (sólo si hacen falta) o en todos los campos siempre

- Carga en SQL

```

1 LOAD DATA [LOW_PRIORITY | CONCURRENT] [LOCAL] INFILE 'file_name'
2   [RELACE | IGNORE]
3   INTO TABLE tbl_name
4   [PARTITION (parctition_name, ...)]
5   [CHARACTER SET charset_name]
6   [{FIELDS | COLUMNS}
7     [TERMINATE BY 'string']
8     [[OPTIONALL] ENCLOSED BY 'char']
9     [ESCAPED BY 'char']
10  ]
11  [LINES
12    [STARTING BY 'string']
13    [TERMINATE BY 'string']
14  ]
15  [IGNORE number {LINES | ROWS}]
16  [(col_name_or_user_var, ...)]
17  [SET col_name = expr, ...]

```

```

1 LOAD DATA LOCAL INFILE "/tmp/Posts.csv"
2   INTO TABLE Posts
3   COLUMNAS TERMINATED BY ',' ENCLOSED BY '"' ESCAPED BY '"'
4   LINES TERMINATED BY '\r\n'
5   IGNORE 1 LINES;

```

- Programación: Lectura con Pandas DataFrame

- Lectura básica de un fichero CSV

```

1 import pandas as pd
2 # Read the CSV file
3 airbnb_data = pd.read_csv("airbnb.csv")

```

- Si queremos establecer la columna `id` como índice

```

1 airbnb_data = pd.read_csv("airbnb.csv", index_col="id")

```

- Si queremos leer solo un conjunto de columnas

```

1 usecols = ["id", "nombre", "barrio", "precio", "noches_minimias"]
2 airbnb_data = pd.read_csv("airbnb.csv", index_col="id", usecols=usecols)

```

- Si queremos indicar un separador de columnas determinado

```

1 usecols = ["id", "nombre", "barrio", "precio", "noches_minimias"]
2 airbnb_data = pd.read_csv("airbnb.csv", index_col="id", usecols=usecols, sep="|")

```


- Si queremos definir el tipo de datos de determinadas columnas

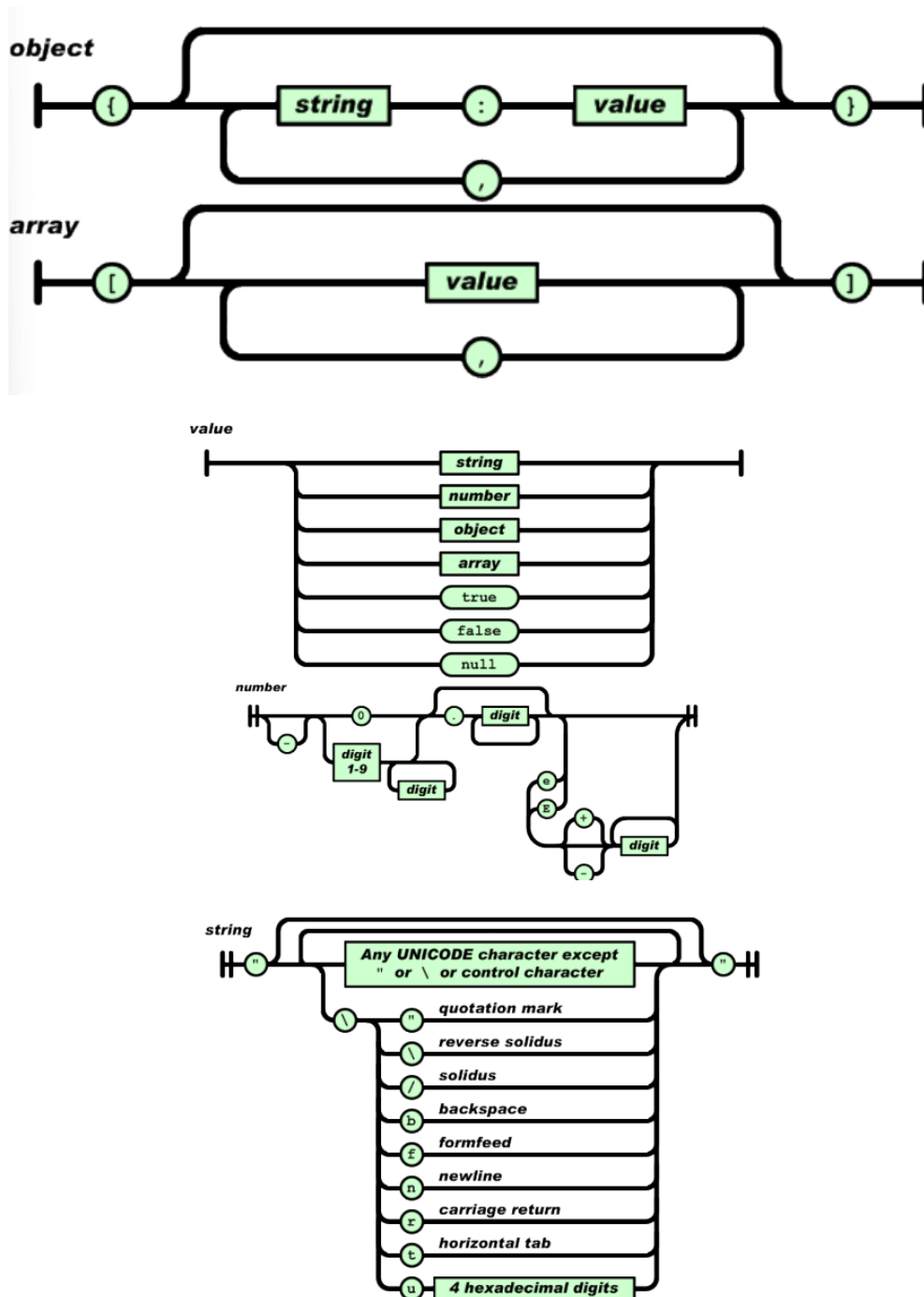
```
1 usecols = ["id", "nombre", "barrio", "precio", "noches_minimias"]
2 airbnb_data = pd.read_csv("airbnb.csv", index_col="id", usecols=usecols, sep="|",
    dtype: {'precio': float, 'barrio': str, 'noches_minimias': int}, decimal=',')
```

1.2.4) JSON (JavaScript Object Notation)

Es, al mismo tiempo, un formato de archivo estándar abierto y un formato de intercambio de datos.

JSON se utiliza a menudo cuando los datos se envían desde un servidor a una página web.

En muchas ocasiones se define a JSON como *autodescriptivo* y fácil de entender.



- JSON Lines (JSONL)

JSON Lines es un formato práctico para almacenar datos estructurados que pueden procesarse de uno en uno.

Ficheros con extensión `.jsonl`

Es un gran formato para archivos de registros.

Sigue las mismas convenciones que JSON salvo que el carácter `\n` se usa como delimitador de líneas.

Cada línea de un fichero `.jsonl` es un JSON válido.

```
1 {"nombre": "Gilbert", "victorias": [["escalera"], ["pareja"]]}
2 {"nombre": "Alexa", "victorias": [["dobles parejas"], ["dobles parejas"]]}
3 {"nombre": "Maya", "victorias": []}
4 {"nombre": "Marisa", "victorias": [["trio"]]}
```

• Lectura con Pandas DataFrame

- Lectura básica de un fichero JSON

```
1 import pandas as pd
2 df = pd.read_json('data.json')
```

- Por defecto, Pandas sigue una orientación basada en columnas en la lectura de los ficheros {columna -> {índice -> valor}}

```
1 json_str = '{"Cursos":{"r1":"Spark"},"Tasas":{"r1":"25000"},"Duracion":{"r1":"50
   días"}}'
2 df = pd.read_json(json_str)
3 print(df)
```

	Cursos	Tasas	Duracion
r1	Spark	25000	50 días

- Otras opciones son `index`, `records`, `split` y `values`.

- Si usamos la orientación `records` ...

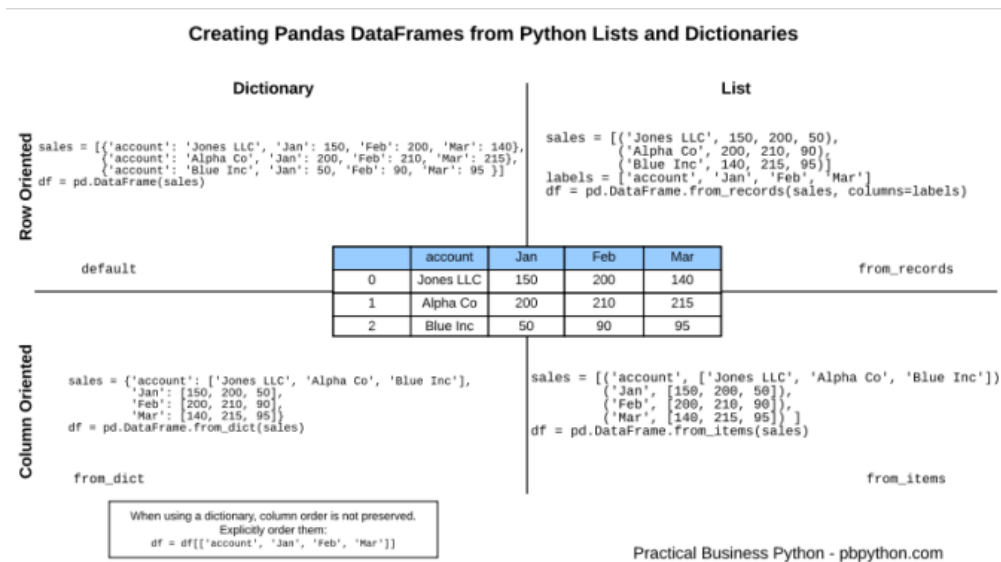
```
1 json_str = '[{"Cursos":"Spark","Tasas":"25000","Duracion":"50
   días","Descuento":"2000"}]'
2 df = pd.read_json(json_str, orient='records')
3 print(df)
```

	Cursos	Tasas	Duracion	Descuento
0	Spark	25000	50 días	2000

- Estableciendo el parámetro `lines` a `True`

```
1 df = pd.read_json('cursos.json', orient='records', nrows=2, lines=True)
```

• Resumen



1.2.5) Apache Avro

Avro ofrece una serie de características:

- Estructuras de datos ricas
- Un formato de datos compacto, rápido y binario
- Un lenguaje de especificación de esquema (Proto + IDL)
- Un formato de fichero contenedor para almacenar datos
- Un esquema de llamada a procedimiento remoto (RPC)
- Integración con lenguajes dinámicos, donde no hace falta recompilar el IDL (opcionalmente se puede hacer para lenguajes estáticos)

1) Schema

- Una declaración de esquema se escribe en JSON. Bien un nombre de un objeto predefinido, o bien un objeto JSON de la forma:

```
1 { "type": "nombreTipo", ... atributos ... }
```

- Ofrece tipos primitivos y complejos. Primitivos: null, boolean, int, long, float, bytes, string, . . .
- `"type": "string"` es equivalente a `"string"`.
- Tipos complejos permitidos: registros (records), enums, mapas, arrays, uniones.
- Tipos complejos: Récor ds (registros)
 - El campo `type` se establece a `"record"`.
 - El campo `name` establece su nombre.
 - El campo `fields` establece sus campos. Cada campo: