# From CSV to RDF

### Here be acronyms...

Julian Padget

# 5 Star Info (Berners-Lee)



   ★ make your stuff available on the Web (whatever format) under an open license

  ★★ make it available as structured data (e.g., Excel instead of image scan of a table)

 ★★★ make it available in a non-proprietary open format (e.g., CSV as well as of Excel)

★★★★ use URIs to denote things, so that people can point at your stuff
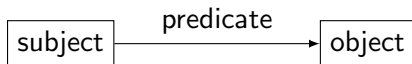
★★★★★ link your data to other data to provide context

# An illustration: DBpedia

- structured, semantically annotated content from Wikipedia
- semantic querying of DBpedia *and other sources*

```
 1 import rdflib
 2 g=rdflib.Graph()
 3 g.load('http://dbpedia.org/resource/Semantic_Web')
 4 semweb=rdflib.URIRef('http://dbpedia.org/resource/Semantic_Web')
 5 dbpedia=rdflib.Namespace('http://dbpedia.org/ontology/')
 6 abstracts=list(x for x in g.objects(semweb, dbpedia['abstract'])
 7    if x.language=='en')
 8 print abstracts[0].value
 9 abstracts=list(x for x in g.objects(semweb, dbpedia['abstract'])
10    if x.language=='ar')
11 print abstracts[0].value
```

# Graphs and Triples

- Basic building block: subject-predicate-object triple

```
              predicate
subject  ──────────────────▶  object
```

- Many syntaxes, including XML
- The fundamental concepts of RDF are:
  - resources: the boxes (more later)
  - properties: the labels on arrows (more later)
  - statements: the s-p-o combination
- Graph is a set of triples
- Relations can be: 1-1, 1-many, many-1, many-many
- Query: triples that satisfy conditions, like SQL
- Resource description with user vocabularies: you define subject, predicate and object

# Resources

- A resource as an object, a "thing" we want to talk about:
  - For example: house, insulation, heating system, ...
- Building block comes from web technology
- Every resource has a URI, a Universal Resource Identifier
- A URI can be:
  - A URL (Web address) or
  - some other kind of unique identifier

# Properties

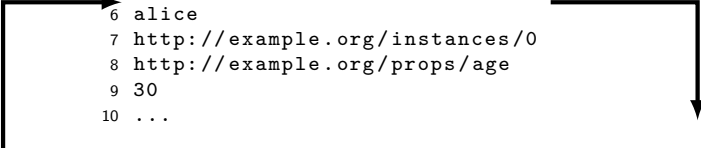- Properties are a special kind of resources
- Used to describe relations between resources
  - For example: "type", "construction date, "storeys", etc.
- Properties also identified by URIs
- Advantages of using URIs:
  - Global, worldwide, unique naming scheme
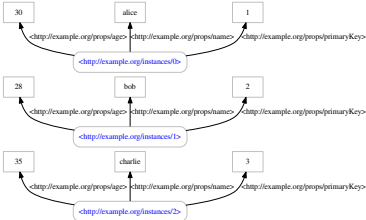  - Limits impact of homonym problem

# Statements

- Statements assert the properties of resources
- A statement is an subject-predicate-object triple comprising:
    - a resource
    - a property, and
    - a value
- Literals are atomic values (strings)
- Values can be resources or literals

# Converting CSV

```
1  http://example.org/instances/0
2  http://example.org/props/primaryKey
3  1
4  http://example.org/instances/0
5  http://example.org/props/name
6  alice
7  http://example.org/instances/0
8  http://example.org/props/age
9  30
10 ...
```

## or as XML

- Convert csv to rdf

```
csv2rdf example.csv > example.rdf
```

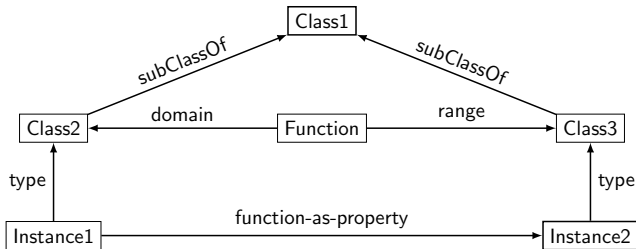- Convert rdf/n3 to rdf/xml

```
1 g=rdflib.Graph()
2 g.load('xls/example.rdf',format='n3')
3 print g.serialize()
```

- Giving

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <rdf:RDF
3    xmlns:ns1="http://example.org/props/"
4    xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
5 >
6   <rdf:Description rdf:about="http://example.org/instances/0">
7     <ns1:name>alice</ns1:name>
8     <ns1:age>30</ns1:age>
9     <ns1:primaryKey>1</ns1:primaryKey>
10   </rdf:Description>
11   <rdf:Description rdf:about="http://example.org/instances/2">
12     <ns1:name>charlie</ns1:name>
13     <ns1:age>35</ns1:age>
14     <ns1:primaryKey>3</ns1:primaryKey>
```

# RDF Schema

- ▶ No assumed or defined domain semantics
- ▶ User-defined by Schema for:
    - ▶ Classes and Properties
    - ▶ Class Hierarchies and Inheritance
    - ▶ Property Hierarchies
- ▶ In practice:
    - ▶ Class, subClassOf, type
    - ▶ Property, subPropertyOf
    - ▶ domain, range



- ▶ Makes simple inference possible...

# RDF(S) Semantics

- The (obvious) inference rules:

$$X \; R \; Y \wedge \mathrm{dom}(R) = T \Rightarrow X \; \mathsf{IsOfType} \; T$$
$$X \; R \; Y \wedge \mathrm{range}(R) = T \Rightarrow Y \; \mathsf{IsOfType} \; T$$
$$T1 \; \mathsf{SubClassOf} \; T2 + T2 \; \mathsf{SubClassOf} \; T3 \Rightarrow T1 \; \mathsf{SubClassOf} \; T3$$
$$X \; \mathsf{IsOfType} \; T1 + T1 \; \mathsf{SubClassOf} \; T2 \Rightarrow X \; \mathsf{IsOfType} \; T2$$

- Given a set of triples can infer other statements
  - Aspirin isOfType Painkiller + Painkiller subClassOf Drug ⇒ Aspirin isOfType Drug
  - Aspirin alleviates Headache + alleviates range Symptom ⇒ Headache isOfType Symptom
- Some triple stores do this automatically, others do not

# Querying a set of triples with SPARQL

```
1 import rdflib
2 g=rdflib.Graph()
3 g.load('xls/example.rdf',format='n3')
4 for row in g.query(
5         'select ?s ?p ?o where { ?s ?p ?o . }'):
6     print row.s
7     print row.p
8     print row.o
```

```
 1 http://example.org/instances/2
 2 http://example.org/props/primaryKey
 3 3
 4 http://example.org/instances/1
 5 http://example.org/props/name
 6 bob
 7 http://example.org/instances/1
 8 http://example.org/props/age
 9 28
10 http://example.org/instances/2
11 http://example.org/props/name
12 charlie
```

# Querying a set of triples with SPARQL

```
1 import rdflib
2 g=rdflib.Graph()
3 g.load('xls/example.rdf',format='n3')
4 ns=dict(props=rdflib.Namespace('http://example.org/props/'))
5 for row in g.query(
6          'select ?s ?o where { ?s props:name ?o .}',initNs=ns):
7     print row.o
```

```
1 charlie
2 alice
3 bob
```

# Triple stores: databases for triples

- A purpose-built database for the storage and retrieval of Resource Description Framework (RDF) metadata
- Implementation: special-purpose (e.g. Jena, JRDF, 4store) or on top of conventional SQL databases
- Advantage: unstructured data—no need to design table structure in advance, so can handle whatever relations are asserted
- Disadvantage: unstructured data—serious impact on performance due to absence of regularity

# Ontology: DIY or re-use?

- DIY example: csv2rdf invented URIs from CSV data
  - `http://example.org/instances/`
  - `http://example.org/props/primaryKey`
  - `http://example.org/props/name`
  - `http://example.org/props/age`
- not very useful... just text labels
- FOAF = Friend-of-a-Friend, W3C-defined ontology
  - Classes: Agent, Organization, Person, ...
  - Properties: account, age, birthday, currentProject, familyName, homepage, ...
- publish ontology... re-use labels
- Writing ontologies: OWL (Web Ontology Language)
  - High-level language for classes and properties
  - Translates to RDF
  - full OWL, OWL-DL and OWL-Lite: computational complexity of reasoning
- Alignment: is my X the same as your X?

# Summary

- URIs are names for resources/properties
- S-P-O triples connect resources with a property
- CSV translates easily into RDF
- RDF semantics provides simple reasoning
- Query set of triples like a database
- Triples stores for lots of triples
- Ontologies for sharing and re-using names (URIs)