# Profiling Linked Open Data

Data Profiling and Data Cleansing

Anja Jentzsch
Information Systems Group (Prof. Dr. Felix Naumann)

# Outline

- Introduction to Linked Data
  - Data Model
  - Data Variety
  - Example Data Set: DBpedia
- Profiling Linked Data
  - Challenges
  - Comparison: Traditional vs Linked Data Profiling
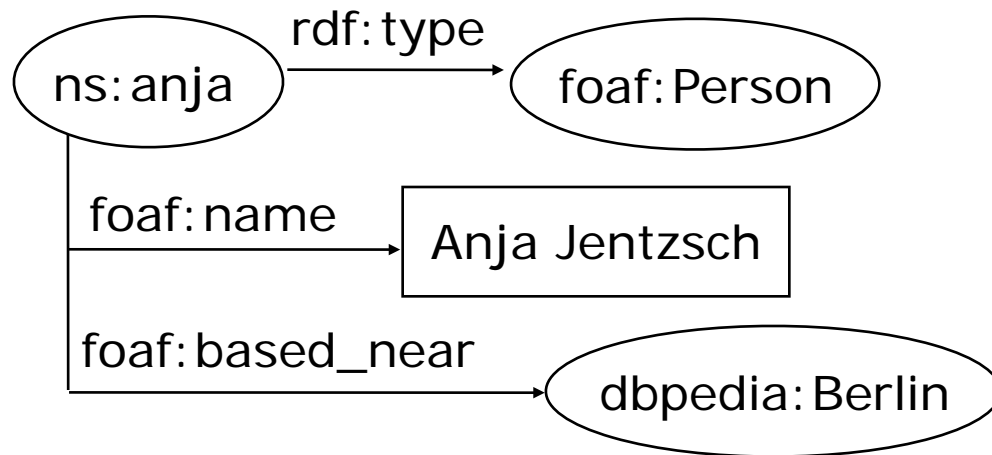  - Existing Approaches

# Linked Data Principles

Set of best practices for publishing structured data on the Web in accordance with the general architecture of the Web.

1.   Use URIs as names for things.
2.   Use HTTP URIs so that people can look up those names.
3.   When someone looks up a URI, provide useful RDF information.
4.   Include RDF statements that link to other URIs so that they can discover related things.

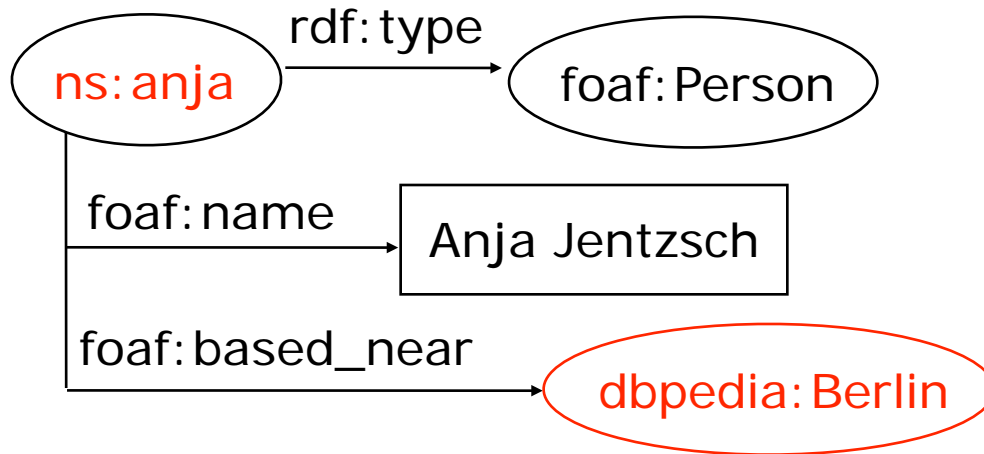Tim Berners-Lee, http://www.w3.org/DesignIssues/LinkedData.html, 2006

ns:anja = http://www.anjeve.de#anja

dbpedia:Berlin = http://dbpedia.org/resource/Berlin

- RDF/XML

```
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-
syntax-ns#"
    xmlns:foaf="http://xmlns.com/foaf/0.1/">
<foaf:Person rdf:about="http://anjeve.de#anja">
    <foaf:name>Anja Jentzsch</foaf:name>
</foaf:Person>
```

- RDF N-Triples
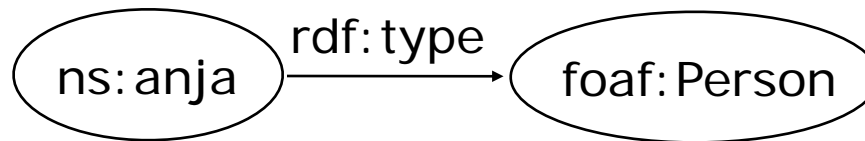
```
<http://anjeve.de#anja> <http://www.w3.org/1999/02/22-
rdf-syntax-ns#type> <http://xmlns.com/foaf/0.1/
Person>   .
<http://anjeve.de#anja> <http://xmlns.com/foaf/0.1/
name> „Anja Jentzsch" .
```
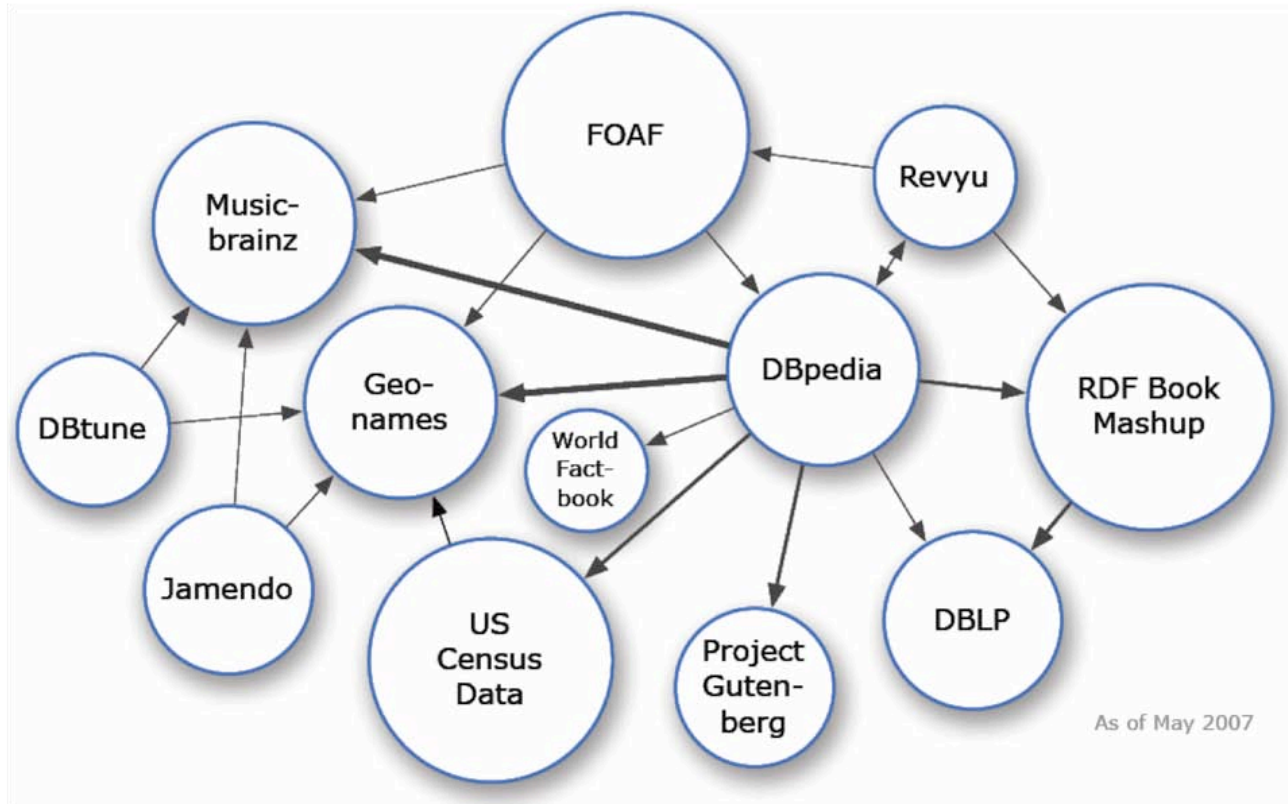
- RDF N3

```
<http://anjeve.de#anja> <http://www.w3.org/1999/02/22-rdf-
syntax-ns#type> <http://xmlns.com/foaf/0.1/Person>  .
```

- \<Subject\> \<Predicate\> \<Object\>
- In the end it's all triples!

# Properties of the Web of Linked Data

- Global, distributed dataspace build on a simple set of standards
  - RDF, URIs, HTTP
- Entities are connected by links
  - Creating a global data graph that spans data sources and
  - Enables the discovery of new data sources
- Provides for data-coexistence
  - Everyone can publish data to the Web of Linked Data
  - Everyone can express their personal view on things
  - Everybody can use the vocabularies/schemas that they like
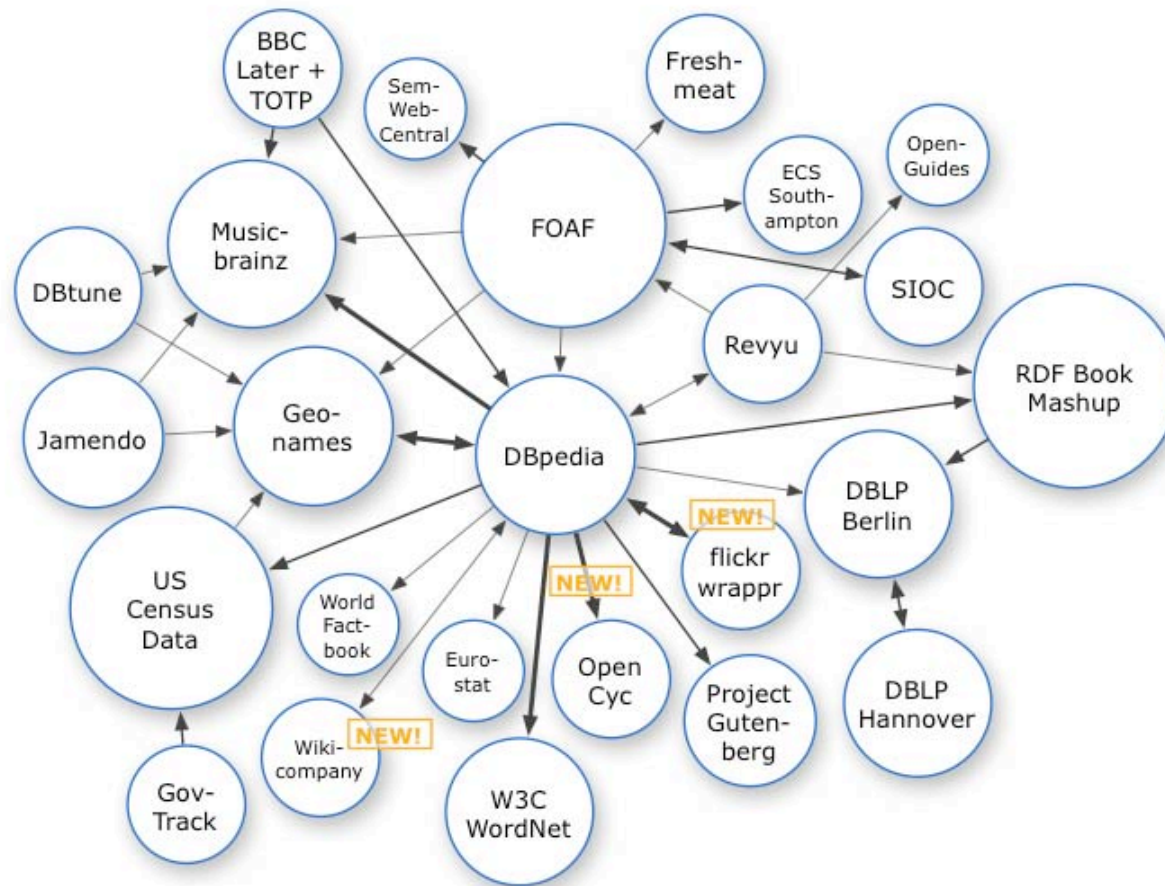
# Web of Data (as of May 2007)



As of May 2007

- 12 data sets
- Over 500 million RDF triples
- Around 120,000 RDF links between data sources
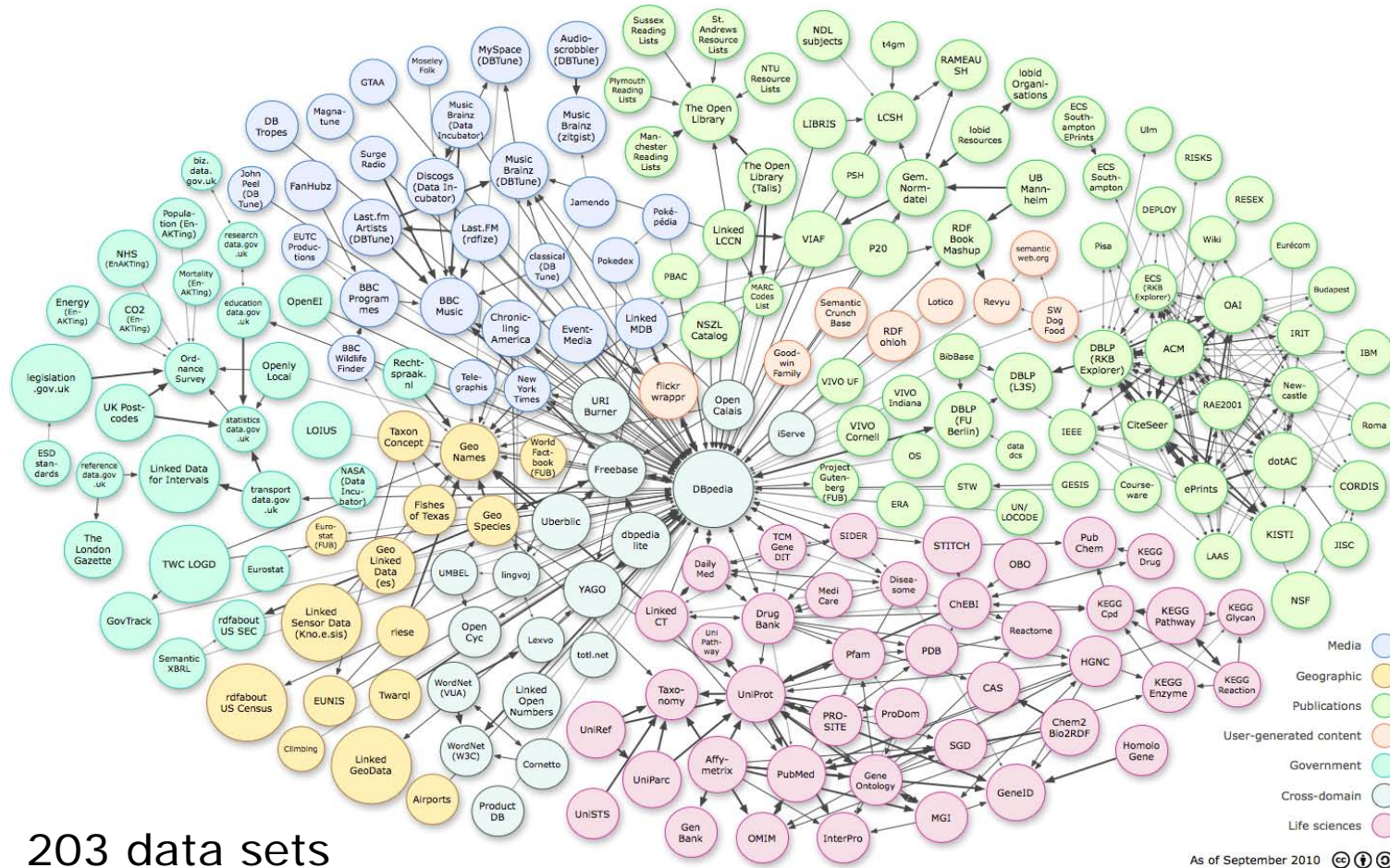
# Web of Data (as of November 2007)



- 28 data sets

13



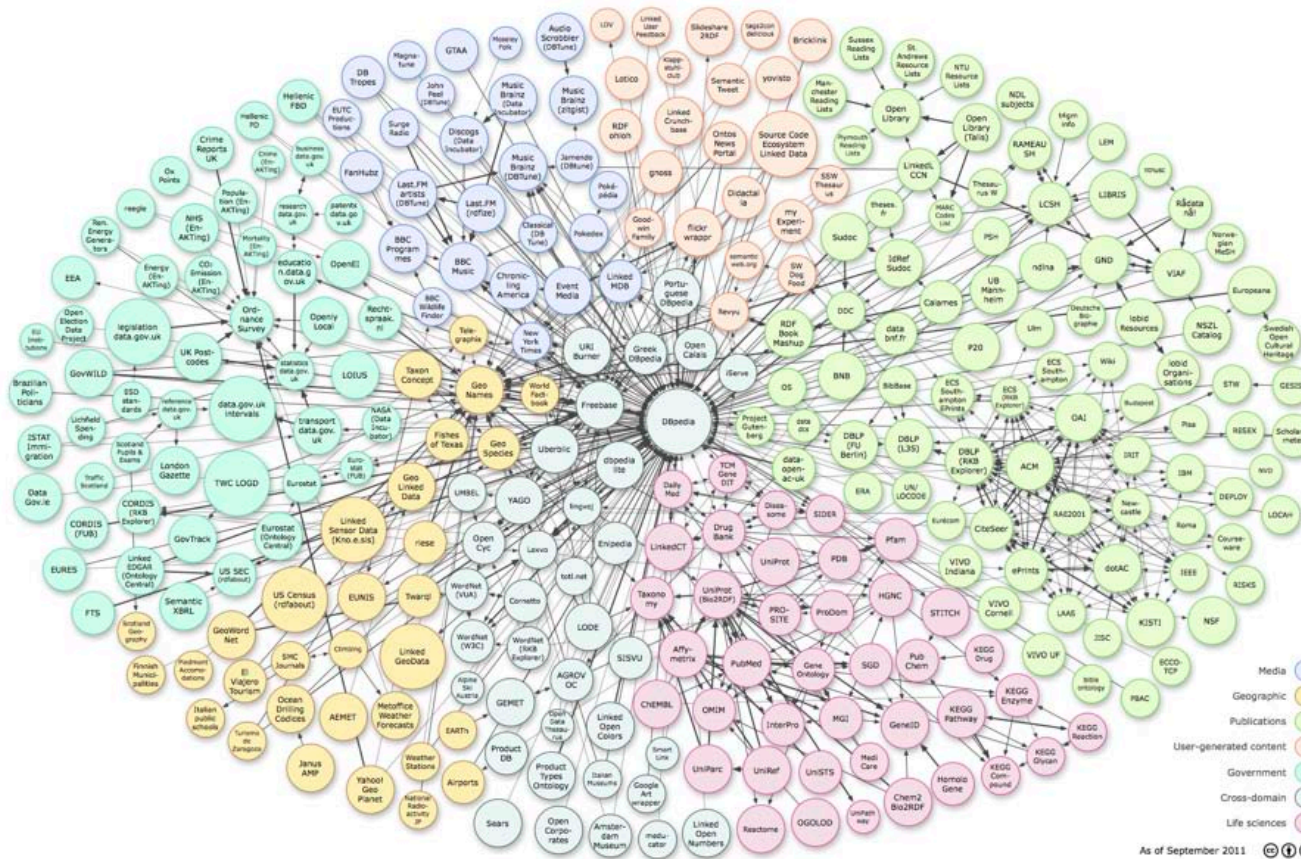As of September 2008

As of July 2009

# Web of Data (as of September 2010)



- 203 data sets
- Over 24,7 billion RDF triples
- Over 436 million RDF links between data sources
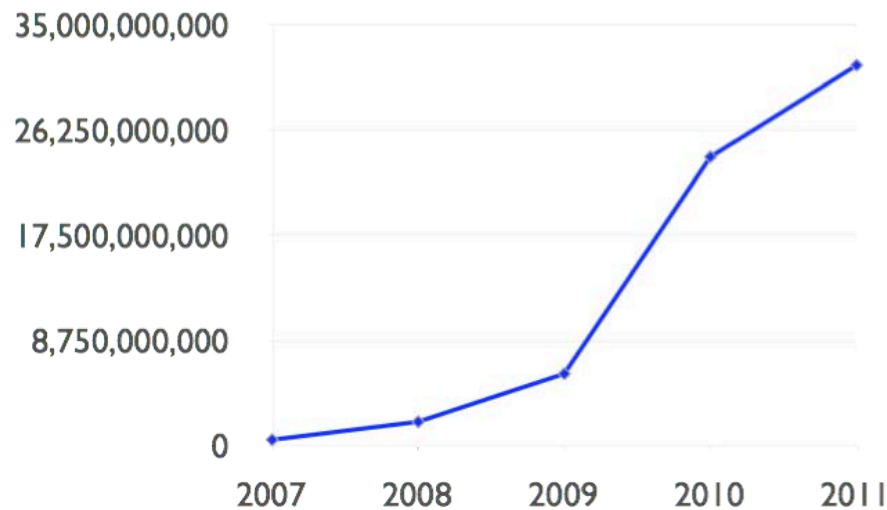
# Web of Data (as of September 2011)



As of September 2011

- 295 data sets

- Over 31 billion RDF triples

- Over 504 million RDF links between data sources

# The Growth in Numbers

| Year | Data Sets | Triples | Growth |
|------|-----------|---------|--------|
| 2007 | 12 | 500,000,000 | |
| 2008 | 45 | 2,000,000,000 | 300% |
| 2009 | 95 | 6,726,000,000 | 236% |
| 2010 | 203 | 26,930,509,703 | 300% |
| 2011 | 295 | 31,634,213,770 | 33% |
| 2013 | ~ 900 | ? | ? |

# Topics on the Web of Data



- LOD Cloud Data Catalog on the Data Hub
  - http://datahub.io/group/lodcloud
- More statistics
  - http://lod-cloud.net/state/

- **The Web of Data is heterogeneous**
  - □ Many different vocabularies are in use (337 as of April 2013)
  - □ Different data formats
  - □ Many different ways to represent the same information

Distribution of the most widely used vocabularies

# Vocabularies on the Web of Data

- Common Vocabularies
  - Friend-of-a-Friend for describing people and their social network
  - SIOC for describing forums and blogs
  - SKOS for representing topic taxonomies
  - Organization Ontology for describing the structure of organizations
  - GoodRelations provides terms for describing products and business entities
  - Music Ontology for describing artists, albums, and performances
  - Review Vocabulary provides terms for representing reviews

21

- Common sources of identifiers (URIs) for real world objects
    - LinkedGeoData and Geonames: locations
    - GeneID and UniProt: life science identifiers
    - DBpedia: wide range of things

# DBpedia - The Hub on the Web of Data

- DBpedia is a joint project with the following goals

  - extracting structured information from Wikipedia

  - publish this information under an open license on the Web

  - setting links to other data sources

- Partners

  - Universität Mannheim (Germany)

  - Universität Leipzig (Germany)

  - OpenLink Software (UK)

UNIVERSITÄT
MANNHEIM

UNIVERSITÄT LEIPZIG

OPENLINK
SOFTWARE

# Extracting structured data from Wikipedia

```
dbpedia:Berlin   rdf:type   dbpedia-owl:City ,

dbpedia-owl:PopulatedPlace ,

dbpedia-owl:Place ;

rdfs:label   "Berlin"@en , "Berlino"@it ;

dbpedia-owl:population   3499879 ;

wgs84:lat   52.500557 ;

wgs84:long   13.398889 .


dbpedia:SoundCloud   dbpedia-owl:location   dbpedia:Berlin .
```

- access to DBpedia data:
  - RDF dumps
  - Linked Data interface
  - SPARQL endpoint

25

1. Improvement of Wikipedia search

2. Data source for applications and mashups

3. Text analysis and annotation

4. Hub for the growing Web of Data

- displays Wikipedia data on map
- aggregates different data sources

# Faceted Wikipedia Search

- Faceted browsing and free text search

## Demo

Berlin is the capital city of Germany and is one of the 16 states of Germany. With a population of 3.45 million people, Berlin is Germany's largest city. It is the second most populous city proper and the seventh most populous urban area in the European Union. Located in northeastern Germany, it is the center of the Berlin-Brandenburg Metropolitan Region, which has 4.4 million residents from over 190 nations. Located in the European Plains, Berlin is influenced by a temperate seasonal climate. Around one third of the city's area is composed of forests, parks, gardens, rivers and lakes.

Back    Confidence: 0.5    Support: 30

Types: Place, Person, Work, Organisation, Species, all other types, untyped

http://spotlight.dbpedia.org

As of September 2011

# The DBpedia Data Set

- Information on more than 3.77 million "things"

  □ 764,000 persons

  □ 192,000 organisations

  □ 573,000 places

  □ 112,000 music albums

  □ 72,000 movies

  □ 202,000 species

- overall more than 1 billion RDF triples

  □ title and abstract in 111 different languages

  □ 8,000,000 links to images

  □ 24,400,000 links to external web pages

  □ 27,200,000 links to other Linked Data sets

# Editing Berlin

```
{{About|the capital of Germany}}
{{Use dmy dates|date=July 2012}}
{{pp-move-indef}}
{{Infobox German state
|Name  =Berlin
|German_name=
|image_photo=Overview Berlin.jpg
|image_caption=Left to right: [[Berliner Fernsehturm]] and Skyline, [[Siegessäule]], [[Kaiser-W
|state_coa =Coat of arms of Berlin.svg
|coa_size =70
|map  =Berlin in Germany and EU.png
|map_size =270
|map_text =Location within [[European Union]] and Germany
|flag  =Flag_of_Berlin.svg
|area  =891.85
|area_source=
|population=3510032{{Verify source|date=August 2012}}
|pop_ref =<ref name="Population">{{cite web|url=http://www.statistik-berlin-brandenburg.de//Pub
Bezirken|work=[[Amt für Statistik Berlin-Brandenburg]]|date=31 October 2011|accessdate=3 March
|pop_date =31 March 2012
|pop_metro =5,963,998
|elevation=34
|demonym =Berliner
|GDP  =94.7
|GDP_year =2010
```

Content that violates any copyrights will be deleted. Encyclopedic content must be **verifiable**.

By clicking the "Save Page" button, you agree to the Terms of Use, and you irrevocably agree to release your contribution under th

Edit summary (Briefly describe the changes you have made)

☐ This is a minor edit (what's this?)  ☑ Watch this page

**Save page**  |  Show preview  |  Show changes  |  Cancel | Editing help (opens in new window)

If you do not want your writing to be edited, used, and redistributed at will, then do not submit it here. All text that you did not write yourself, except brief excerp

---

### Berlin
— State of Germany —

Coordinates: 52°30'2"N 13°23'56"E

| | |
|---|---|
| Country | Germany |
| **Government** | |
| - Governing Mayor | Klaus Wowereit (SPD) |
| - Governing parties | SPD / The Left |
| - Votes in Bundesrat | 4 (of 69) |
| **Area** | |
| - City | 891.85 km² (344.3 sq mi) |
| Elevation | 34 - 115 m (-343 ft) |
| **Population** (30 April 2011)[1] | |
| - City | 3,471,756 |
| - Density | 3,892.8/km² (10,082.2/sq mi) |
| - Metro | 4,429,847 |

- Companies in DBpedia

- Def. 1: Subject having a predicate `companyName`
  - → 14,292

- Def. 2: Subject having a category that starts with `'compan'`
  - → 21,753

- Def. 3: Subject having a `wikiPageUsesTemplate` with value `Template:infobox_company`
  - → 15,491



1.companyName
222

759    3,207

10,104

9,204    494

1,686

2.compan%
category

3.company
template

34

- **DBpedia**: `?c wikiPageUsesTemplate Template:infobox_company`


- 1,083 different attributes
- 499 appear only once



- 39 distinct ones contain `name` as substring

  `companyName, commonName, publicName, …`

- 273 companies without any name attribute

# DBpedia Company Attribute Distribution



| | | | |
|---|---:|---|---:|
| location | 20617 | companyName | 13355 |
| products | 18176 | name | 2036 |
| wikiPageUsesTemplate | 18048 | surname | 25 |
| keyPeople | 17836 | railroadName | 8 |
| industry | 16822 | companyNickname | 4 |
| foundation | 15826 | pastNames | 4 |
| homepage | 14476 | absNameProperty | 3 |
| companyType | 13433 | dnvNameProperty | 3 |
| companyName | 13355 | labelName | 3 |
| companyLogo | 9006 | logoFilename | 3 |
| numEmployees | 6207 | dvdEuroCompanyName | 2 |
| revenue | 5030 | filename | 2 |
| locationCity | 4098 | longName | 2 |
| locationCountry | 3212 | websitename | 2 |
| companySlogan | 2815 | alternativeNames | 1 |
| areaServed | 2557 | birthname | 1 |
| relatedInstance | 2284 | brandName | 1 |
| type | 2152 | bTcgvuvCompanyName | 1 |
| parent | 2054 | companyNameLocal | 1 |
| name | 2036 | companyNamesBigBum | 1 |
| netIncome | 1663 | europeanTradeAssociationCompanyName | 1 |
| founder | 1597 | familyCorporationCompanyName | 1 |
| subsid | 1232 | formerNames | 1 |
| nihongoProperty | 1141 | fukCompanyName | 1 |
| slogan | 1087 | golfFacilityName | 1 |
| coorTitleDmsProperty | 960 | hangulName | 1 |
| logo | 925 | iceCreamCompanyName | 1 |
| services | 904 | nativeName | 1 |
| operatingIncome | 896 | nickname | 1 |
| owner | 680 | officialName | 1 |
| otheruses4Property | 510 | oldName | 1 |
| intl | 503 | organisationName | 1 |
| forProperty | 467 | publicCompanyName | 1 |
| divisions | 429 | renamed | 1 |
| date | 422 | shortName | 1 |
| locations | 419 | wineryName | 1 |

- since March 2010 collaborative editing of
  - DBpedia ontology
  - mappings from Wikipedia infoboxes and tables to DBpedia ontology
- curated in a public wiki with instant validation methods
  - http://mappings.dbpedia.org
- multi-langual mappings to the DBpedia ontology:
  - ar, bg, bn, ca, cs, de, el, en, es, et, eu, fr, ga, hi, hr, hu, it, ja, ko, nl, pl, pt, ru, sl, tr

- allows for a significant increase of the extracted data's quality
  - each domain has its experts

# DBpedia Ontology

- 359 classes

  - 2,347 mappings from Wikipedia infoboxes to ontology classes (overall)

- 800 object properties, 859 datatype properties, 116 specialized datatype properties

  - 5,859 mappings from Wikipedia infobox properties to ontology properties (en)

- 45 owl:equivalentClass and 31 owl:equivalentProperty mappings to http://schema.org

- Example: Wikipedia/DBpedia
- Schema chaos: Many attribute synonyms
  - Hundreds of different attributes
  - `companyName` vs. `organizationName` vs. `name` vs. `company`
- Schema misuse: Many attribute homonyms
  - `foundation` attribute in DBPedia may contain
    - ◇ Person who founded the company
    - ◇ Year/Date company was founded
    - ◇ Location where the company was found

- Linked Data published by third parties
    - Personal view on data
    - Misinterpretation

- Loosely defined schema
    - Missing property definitions
    - Property types used inconsistently

# Outline

- Introduction to Linked Data
  - Data Model
  - Data Variety
  - Example Data Set: DBpedia
- **Profiling Linked Data**
  - Challenges
  - Comparison: Traditional vs Linked Data Profiling
  - Existing Approaches

41

- Current situation:
  - Web of Data is growing

- Advantages:
  - Wealth of information
  - Easy, public access
  - Interesting domains

- Challenges:
    - Heterogeneity
        - ◇ Loose structure    Things have different predicate stets
        - ◇ Incomplete    Subjects do not have name predicate
        - ◇ Poorly formatted    Predicate values have many patterns
        - ◇ Inconsistent    Multiple representations claim opposite
    - Volume of data

- Linked Data integration
- Linked Data publication
- Interlinking Linked Data sets

- Data profiling allows for analyzing
  - Semantic heterogeneity
  - Structural heterogeneity

As of September 2011

- Required knowledge for describing Linked Data sets:
  - Detailed characteristics of a data set (or parts of it)
  - Relevance of data set
  - Retrieving and processing these information for a large number of data sets is practically unfeasable
- Easy finding approach:
  - Popular data sets (e.g. DBpedia, Geonames)
  - Not always optimal:
    - ◇ If data domain is highly specialized and not covered by popular data sets in sufficient detail
    - ◇ If different parts of the data sets are covered by several external data sets (e.g. publications both on computer science (DBLP) and medicine (PubMed))

# Profiling Linked Data - Motivation

- Evolving Linked Data sets require constant re-analysis
- Interlinking Linked Data sets
  - Link discovery problem has been addressed by several approaches (Silk, LIMES, KnoFuss)
  - Published data sets often interlinked with the help of researchers interested in the Linked Data initiative
- Identifying relevant sources did not acquire much attention
- Gathering linkage/integration possibilities is a time-consuming effort
- Reduce effort to perform exploratory search
- Bringing publication and interlinking process closer together

48

- Topic(s)
- Statistical characteristics
  - □ Classes
  - □ Instances
  - □ Properties
  - □ Property values (and distribution)
- Language(s)
- Schema
- Data set granularity
- Relevance
- …

49

- Documents on the data set (website, papers, …)
- Metadata files (VoID / Semantic Sitemap)
- Data registries (The Data Hub)

- ‣ Provide valuable but usually not fine-grained information on content of Linked Data sets

- State of the art Data profiling
    - Based on columns
    - Assumes well-defined semantics
    - Expects regular data

- Heterogeneity on the Web of Data
    - Diverse sources
    - ➡ Diverse structures
    - ➡ Diverse views
- RDF: nested graphs

- Nevertheless some "clean" LOD sources exist (ontologies, RDFS)
- Integration problem remains

- Number of Triples
- Number of Instances
- Average number of properties per instance

# Data Set Statistics: Schema-Based

- Number of classes
- Number of instances per class
- Average number of values per property
- Percentage of top-k properties per class
- Number of different datatypes and language tags used
- Average length of strings (per property)
- Value ranges for numeric properties (per property)
- Ratio URIs/literals as objects
- Co-occuring classes
- Co-occuring properties
- Equivalent classes
- Equivalent properties

53

- Number of different properties per data set and class
- Number of RDF links set between instances of the data set
- Number of RDF links pointing at instances within the dataset
- Number of RDF links pointing at instances in other data setsAverage indegree/outdegree
- Number of links likely pointing at HTML pages

- Number of classes/properties that are reused from common vocabularies
- Percentage of classes/properties that are reused from common vocabularies
- Topic (VoID, Semantic Sitemaps, The Data Hub, …)

# Existing Linked Data Profiling Approaches

- ProLOD

- Creating voiD descriptions

- Finding relevant link target

- Schema induction (gold-miner)

# ProLOD

- Christoph Böhm, Felix Naumann et. al. @ NTII2010, ICDE2010

- Offers profiling methods to deal with loosely structured, unclean and inconsistent data on the Web of Data

- Well-known profiling techniques

- Web-based tool

# ProLOD

- Suite of methods ranging from:
    - Domain level (clustering, labeling)
    - Schema level (matching, disambiguation)
    - Data level (data type detection, pattern detection, value distribution)

- Heterogeneity

- Consider a `height` predicate
    - Average value is 30 (Feet? Inches?)
    - But there are heights of buildings (in feet) and plants (in inches)
    - Average height of a building is 64 feet
    - Average height of a plant is 4 inches

- Prerequisite for meaningful profiling
- Volume of the data

- **Similarity of data entities**
  - □ Schema Similarity = Jaccard Similarity
- **Dissimilarity of data entities**
  - □ Schema Dissimilarity = 1 – Schema Similarity
- **Intra-Cluster Dissimilarity**
  - □ Average pairwise Schema Dissimilarity
- **Cluster Centroid**
  - □ Schema of a cluster = Mean Schema
  - □ Threshold Mean Schema
    = Predicates required to be in t% of subjects
  - □ Top N Mean Schema (default)
    = N most frequent properties (N avg number of properties)

60

- Iterative
  - □ Cluster data with k=2
  - □ While Cluster dissimilarity > threshold
    - ◇ Choose single Cluster C
    - ◇ Cluster C with k=2 (overall k increases)
- Hierarchical
  - □ recursive call of iterative K-Means
  - □ Predefined set of parameters to stop recursion
    - ◇ Max depth: 3
    - ◇ Max number of clusters in depth d: d=0:50, d=1:15, d=2:7
    - ◇ Max Cluster Dissimilarity: 0.3
    - ◇ Min Cluster Size: 100

- Use of textual subject descriptions
  - `rdf:comment`
  - `rdf:about`
  - `shortAbstract` (in DBpedia)

- Top k tf-idf weighted terms (default k=3, cluster is a document)
- Evaluation:
  - Given a grouping by `wikiUsesTemplate`
  - >56% of labels contain token from template name
  - More textual descriptions per cluster → higher percentage

- Top k predicates from Mean Schema

- Enables initial understanding of the actual structure of the data (set of triples does not expose much structural information)

- Determining the actual schema (e.g., distinct attributes of a cluster)

- Finding equivalent attributes (e.g., name, family name, and surname)

- Discovering poor attributes (i.e., those that do not contain useful values for most data entries)

- Discover attribute correlations
  - association rules
  - inverse relations
  - foreign key relationships

63

- Heterogeneity

- Determine set of attributes with 'clean' semantics from initial predicates
- Example: media cluster where entities have different predicates Consider author and/or developer predicates

- Most entities have author and developer, *distinct* semantics
  → Data ok, Clustering ok
- Most entities have either author or developer, *distinct* semantics
  → Data ok, Clustering questionable
- Most entities have author and/or developer, *similar* semantics
  → Data dirty, Clustering ok

- Apriori Algorithm, Agrawal and Srikant, 1994

  - media cluster example:

| Rule | Confidence | Correlation Coefficient |
|---|---|---|
| $genre, isbn \Rightarrow author$ | 0.99 | 0.67 |
| $isbn \Rightarrow author$ | 0.92 | 0.66 |
| $isbn \Rightarrow author, genre$ | 0.83 | 0.66 |
| $author, genre \Rightarrow isbn$ | 0.70 | 0.66 |
| $author \Rightarrow isbn$ | 0.64 | 0.66 |
| $author \Rightarrow genre, isbn$ | 0.58 | 0.67 |

- Conclusion:

  - genre, isbn, author together form part of an entity's schema

  ➡ Assumption: complement each other

  - distinct semantics

- Use of Correlation Coefficient, Antonie and Zaiane, 2004

- media cluster example:
  - name -> not( title )

- Conclusion:
  - Subjects from different domains in cluster → poorly built
    - ◇ Perform (sub)clustering with ProLOD
  - Semantic equivalence of predicates
    - ◇ Merge predicates in ProLOD

66

- Subject X holds link to Subject Y via predicate $X \overset{A}{\to} Y$
- $X \overset{A}{\to} Y$    a$Y \overset{B}{\to} X$   , then A and B are inverse links.

- Example:

| $\underrightarrow{PredicateA}$ | $\underleftarrow{PredicateB}$ | Corr Coef | Frequency |
|---|---|---|---|
| before | after | 0.239 | 28856 |
| sisterStations | sisterStations | 0.749 | 7494 |
| precededBy | followedBy | 0.830 | 7097 |
| spouse | spouse | 0.322 | 1964 |
| before | before | -0.003 | 738 |
| star | exoplanet | 0.895 | 188 |

- Conclusion:
  - Redundancy of e.g. before/after and sisterStations
    - ◇ Fuse with ProLOD
  - Misuse of before
    - ◇ Exclude before with ProLOD

- (mostly) State-of-the-art Profiling for attribute values
- Distinction of values: literals, internal and external links

- Profiling for external links and literals
  - Data types
    (String, Text, Integer, Decimal, Date)
  - String → determine (normalized) patterns
  - Integers, Decimals →  display value ranges
  - Set of user-defined keywords, and context rules
    - ◇ Months: Jan, Feb, Mar …
      - – Markus vs. Mar-06-2010 Aaaaaa vs. MONTH-99-9999
    - ◇ File extensions: .jpg, .mpg, …
    - ◇ URL Schemas: http, ftp, …

68

Clustering
+
Modifications
(merge, split …)

Schema Discovery
 +
Modifications
(filter, fuse, rename …)



Data
+
Views

Profiling

continuous process

Understanding

Metadata

69

http://youtu.be/_qyhVMOTbm0

- Christoph Böhm, Johannes Lorey, Felix Naumann @ ISWC2010

- Scalable approach for segmenting, annotating, and enriching Linked Data sets
- Extend scope of voiD (Vocabulary of Interlinked Datasets)
    - Connected sets
        - 2 resources reside within the same connected dataset, iff there is a link of a specific type between them
    - Conceptual sets
        - 2 resources are contained in the same conceptual dataset, iff they are of the same or of similar type

✓ **void:datset**

✓ **void:linkset**

✓ **void:uriLookupEndpoint**

- based on URI patterns of dataset resources

✓ **dcterms:description**

- based on ranked list of subject types (rdf:type)

✓ **void:exampleResource**

- based on dataset entity providing most statements

✓ **void:statItem**

- various statistical information about dataset

✓ **void:vocabulary**

- based on URIs of predicates

# Connected Datasets for voiD?

# Conceptual Datasets for voiD?



**city dataset**
**capital dataset**

- Andriy Nikolov, Mathieu d'Aquin @ LDOW2011, WWW2011

- Two step approach:
  - Use subset of labels for keyword-based search on Semantic Web indexes to retrieve potentially relevant instances in external data sets
  - Use ontology matching techniques to filter out irrelevant sources by measuring semantic similarities between classes

- Keyword-based search for relevant instances:
  - Randomly select subset of individuals of belonging to a class (reduces number of search queries)
  - Query search engine (Sig.ma) for labels of each instance in subset
    - ◇ Sig.ma returns RDF document with references to instances, their sources and the classes they belong to
  - Aggregate search result
    - ◇ Load Sig.ma RDF documents in store and group instances by their sources
  - Data sets are ranked according to the numbers of returned instances

- Use ontology matching techniques to filter out irrelevant results
  - Use ontology matching algorithm (CIDER) to measure similarity between classes in original data sets and found classes
  - Filter out classes with low similarity index by applying a filter
  - Apply instance-based matching to BTC data set to map schemata based on ow:sameAs relations
  - Merge remaining classes with the classes obtained from the BTC schema mappings
  - Filter only instances that belong to the resulting class set

# gold-miner

- Johanna Völkel, Mathias Niepert. http://code.google.com/p/gold-miner/

- Statistical schema induction
- Steps
  - Terminology acquisition from data set(s): classes and properties
  - Association rule mining
  - Ontology construction

# Other existing approaches

- Conditional inclusion dependencies (Bauckmann, Naumann)
  - DBpedia person analysis in English and German DBpedia
  - Conditions on which German persons occur in English DBpedia

- Web of Data is growing
- Advantages:
  - □ Wealth of information
  - □ Easy, public access
  - □ Interesting domains
- Challenges:
  - □ Heterogeneity
    - ◇ Loose structure
    - ◇ Incomplete
    - ◇ Poorly formatted
    - ◇ Inconsistent
  - □ Volume of data

81

- Christoph Böhm, Felix Naumann, Ziawasch Abedjan, Dandy Fenz, Toni Grütze, Daniel Hefenbrock, Matthias Pohl, David Sonnabend. Profiling Linked Open Data with ProLOD. NTII2010, ICDE2010, 2010.

- Andriy Nikolov , Mathieu D'Aquin. Identifying Relevant Sources for Data Linking using a Semantic Web Index. LDOW2011, WWW2011, 2011.

- Andriy Nikolov, Enrico Motta. Capturing Emerging Relations between Schema Ontologies on the Web of Data. COLD2010, ISWC2010, 2010.

- Johanna Völkel, Mathias Niepert. Statistical Schema Induction. ESWC2011, 2011.