

Max Planck Institute
for Human Development

From relational databases to linked data:R for the semantic web

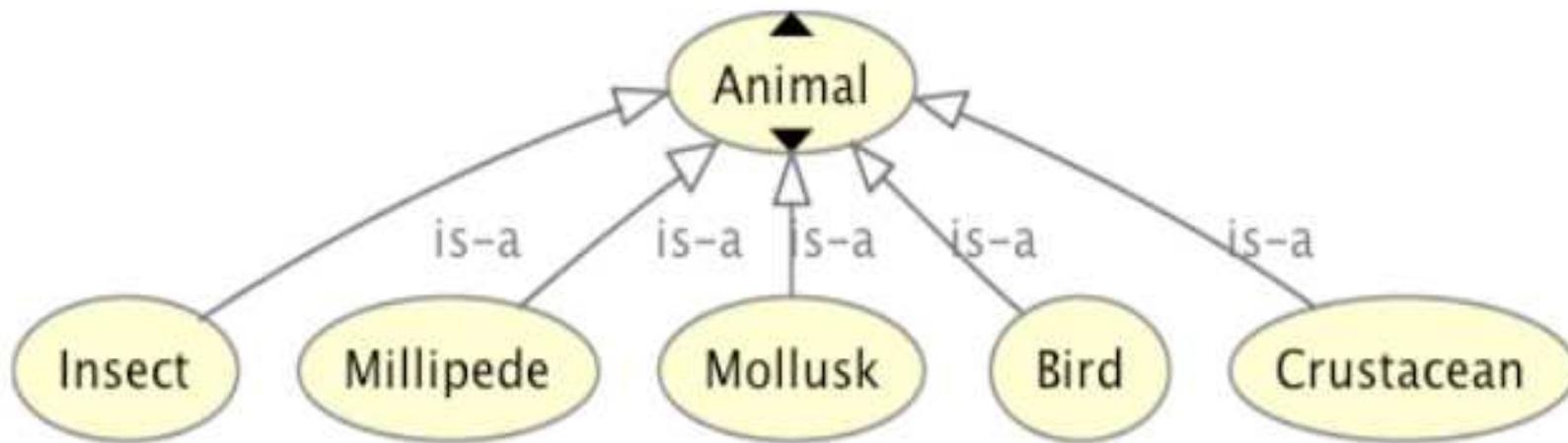
Jose Quesada,
Max Planck Institute, Berlin



Who this talk targets

- You have big data; you use a database
- You have an evolving schema definition. Sometimes at runtime
- You are interested in alternative ways to present your data
- You would thrive by using data out there, if only they were more accessible

Semantic web



Credit: Jim Hendler

THE TWO TOWERS



The Semantic web



- Ontology as Barad-dur (Sauron's tower)
 - Extremely powerful
 - Patrolled by Orcs
 - Let one little hobbit in it, and the whole thing could come crashing down
 - OWL

The Semantic web



- Ontology as Barad-dur (Sauron's tower)

- Extremely powerful

Decidable logic basis

- ~~Patrolled by Orcs~~

- Let one little ~~hobbit~~^{inconsistency} in it, and the whole thing could come crashing down

- OWL

Inconsistency

File View Bookmarks Resource Holder Advanced About

Address: http://swccf.jpl.nasa.gov/ontology/process.owl

OWL Ontology: [process.owl](#)

Annotations:
owl:versionInfo : 1.0

Total Number of Classes: 1537 (Defined: 1537, Imported: 0)
Total Number of Datatype Properties: 19 (Defined: 19, Imported: 0)
Total Number of Object Properties: 102 (Defined: 102, Imported: 0)
Total Number of Annotation Properties: 2 (Defined: 2, Imported: 0)
Total Number of Individuals: 150 (Defined: 150, Imported: 0)

Advanced Ontology Statistics:

| General Statistics | Property Tree Statistics | Satisfiable Class-Tree Statistics |
|---|--------------------------|-----------------------------------|
| No. of Unsatisfiable Classes: 1 | | |
| DL Expressivity: ALCHOFDL | | |
| No. of GCs: 2 | | |
| No. of Sub-classes: 1928 | | |
| No. of Disjoint Axioms: 1 | | |
| No. of Functional Properties: 10 | | |
| No. of Inverse Functional Properties: 0 | | |
| No. of Transitive Properties: 0 | | |
| No. of Symmetric Properties: 0 | | |
| No. of Reflexive Properties: 0 | | |

Axioms causing the inference

- 1) OceanCrustLayer \sqsubseteq owl:Nothing
- 2) \perp (OceanRegion \sqsubseteq TopographicRegion)
- 3) \perp (TopographicRegion \sqsubseteq EarthRegion)
- 4) \perp (EarthRegion \sqsubseteq Region)
- 5) \perp (Region \sqsubseteq GeometricalObject_2D)
- 6) \perp (GeometricalObject_2D \sqsubseteq (HasDimension . (**)^n <xsd:integer>))
- 7) (OceanCrustLayer \sqsubseteq CrustLayer)
- 8) \perp (CrustLayer \sqsubseteq LithospheresLayer)
- 9) \perp (LithospheresLayer \sqsubseteq SolidEarthLayer)
- 10) \perp (SolidEarthLayer \sqsubseteq Layer)
- 11) \perp (Layer \sqsubseteq GeometricalObject_3D)
- 12) \perp (GeometricalObject_3D \sqsubseteq (HasDimension . (**)^n <xsd:integer>))

Hide out relevant parts of axioms

owl:Thing
Air
Aquifer
ArrayOrNill
ChemicalSubstance
CoastalRegion
Crop
Cryosol
DepthHoar
DewPoint
DustOrSendOrSpray
EarthRealm
EcologicalProcess
[ElectromagneticRad
EquipmentCharacter
Evapotranspiration
FieldStrength
Fishery
FormOfSubstanceOr
FreezingRain
GeoReferenceInform
Glacier
Grid
Ground_substance
Hazard
HazardousLevel
[Horizon_Profile]

The semantic web



- The tower of Babel
 - We will build a tower to reach the sky
 - We only need a little ontological agreement
 - Who cares if we all speak different languages?

This is RDFS

Statistics matter here

Web-scale

Lots of data; finding anything in the mess can be a win

Approaches to data representation

- Objects
- Tables (relational databases)
- Non-relational databases
- Tables (data.frame)
- Graphs

What one can do with semantic web data, now:

People that died in Nazi Germany and if possible,
any notable works that they might have created

```
SELECT *  
WHERE {  
  ?subject dbpprop:deathPlace  
<http://dbpedia.org/resource/Nazi_Germany> .  
  OPTIONAL {  
    ?subject dbpedia-owl:notableworks ?works  
  }  
}
```

| subject | works |
|---|---|
| <u>:Anne Frank</u> | <u>:The Diary of a Young Girl</u> |
| <u>:Martin Bormann</u> | - |
| <u>:Irene Neimirovsky</u> | - |
| <u>:Erich Fellgiebel</u> | - |
| <u>:Friedrich Ferdinand%2C Duke of Schleswig-Holstein</u> | - |
| <u>:Friedrich Olbricht</u> | - |
| <u>:Ludwig Beck</u> | - |
| <u>:Erwin Rommel</u> | - |
| <u>:Maurice Bavaud</u> | - |
| <u>:Early Years of Adolf Hitler</u> | - |
| <u>:Emil Zegad%82owicz</u> | - |
| <u>:Friedrich Fromm</u> | - |
| <u>:Heimuth James Graf von Moltk</u> | |



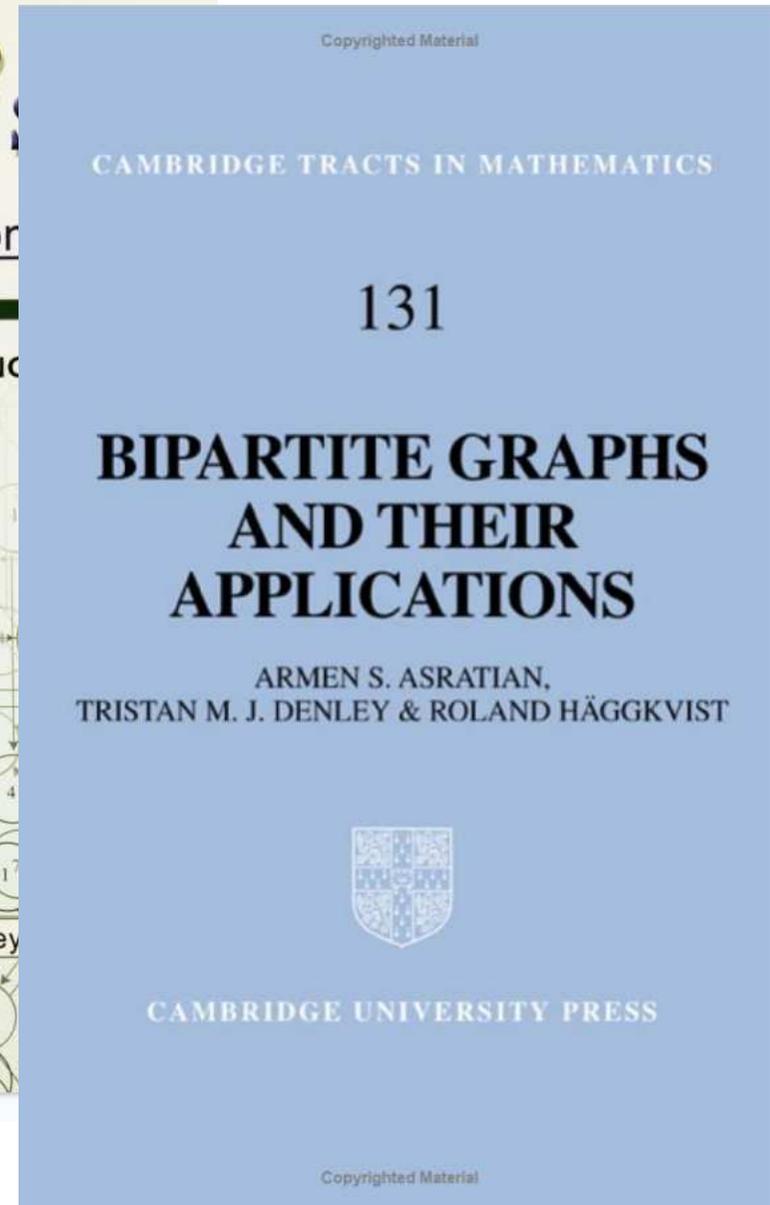
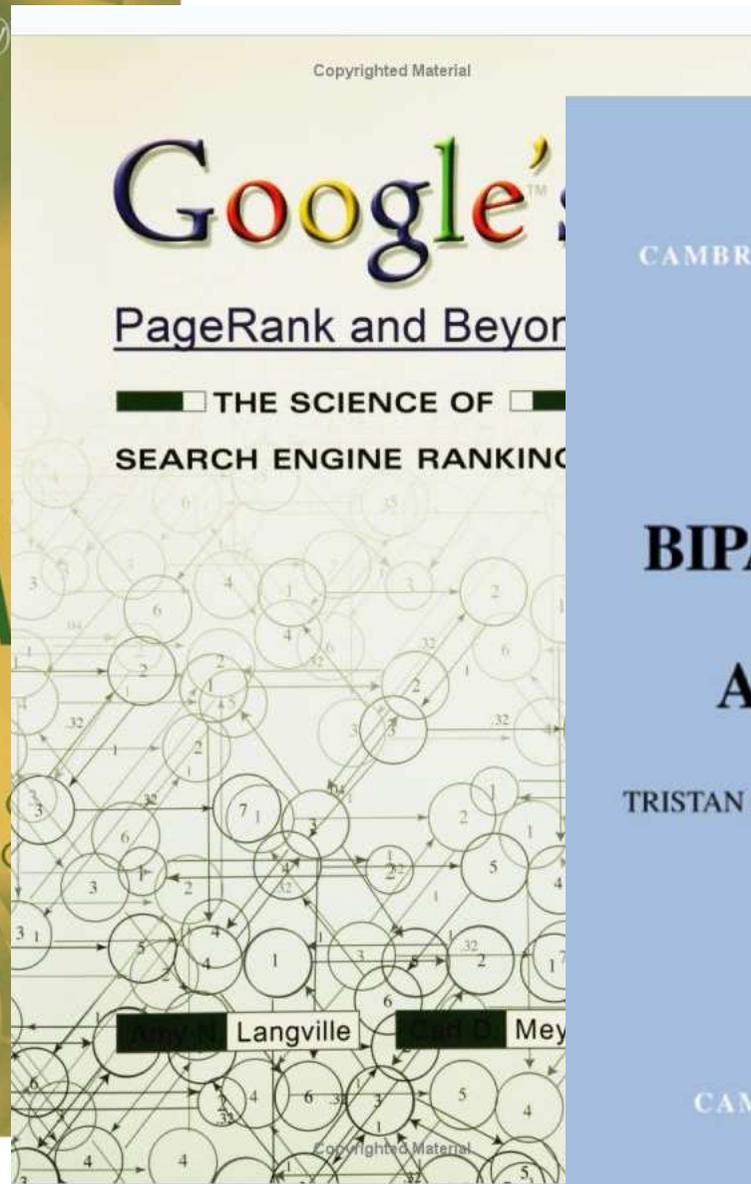
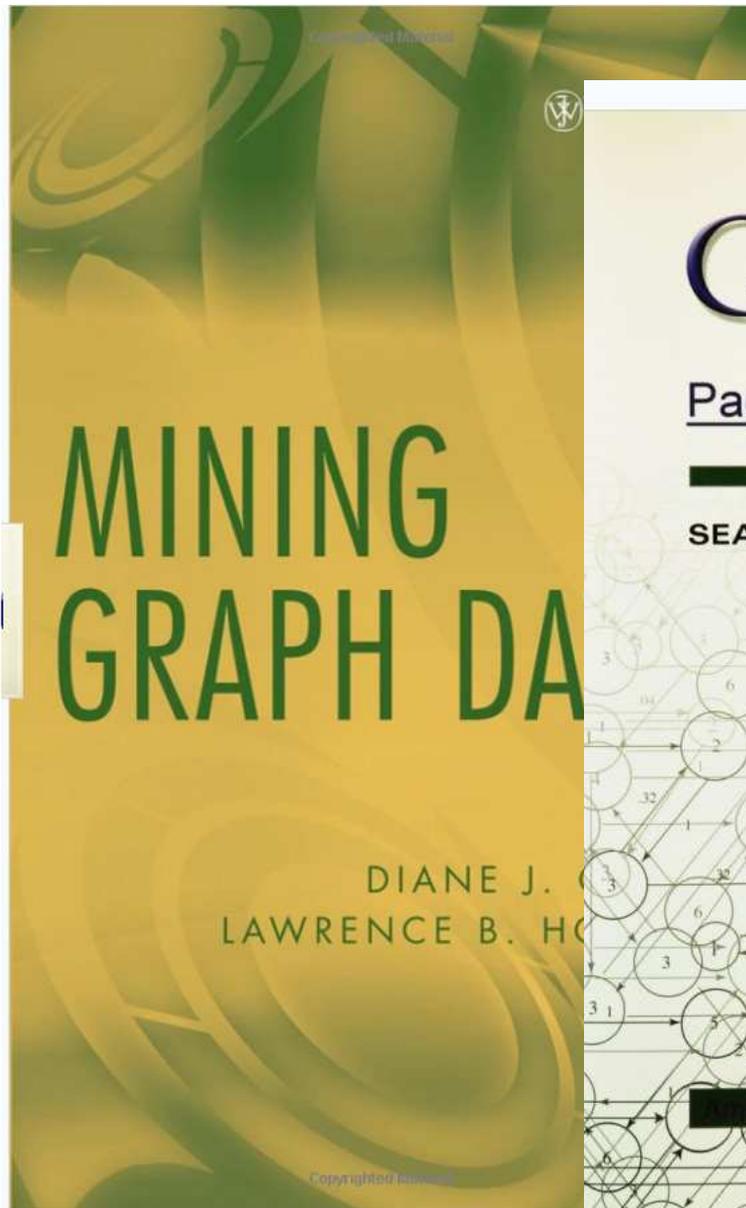
- Scale to the entire web
- Do reasoning with open word assumption
- Retrieval in real-time
- Go beyond logics
- Use cases:
 - Real time city
 - Cancer monographs for WHO
 - Gene expression finding

RDF is a graph

- We have lots of interesting statistics that run on graphs
- In many Semantic Web (SW) domains a tremendous amount of statements (expressed as triples) might be true but, in a given domain, **only a small number of statements** is known to be true or **can be inferred to be true**. It thus makes sense to attempt to estimate the truth values of statements by **exploring regularities** in the SW data with machine learning

Scale

- You cannot use the entire thing at once:
subsetting
- Are there patterns in knowledge structures
that we can use for subsetting?



Idea

- Graph theory applied to subsetting large graphs
- Developing Semantic Web applications requires handling the RDF data model in a programming language
- Problem: current software is developed in the object-oriented paradigm, programming in RDF is currently triple-based.



Data

IMDB is a big graph:

- 1.4 m movies
- 1.7 m actors
- 11 M connections
 - Movies have votes
- Bipartite network

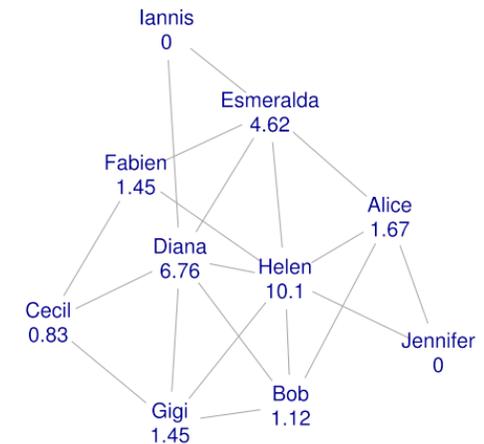
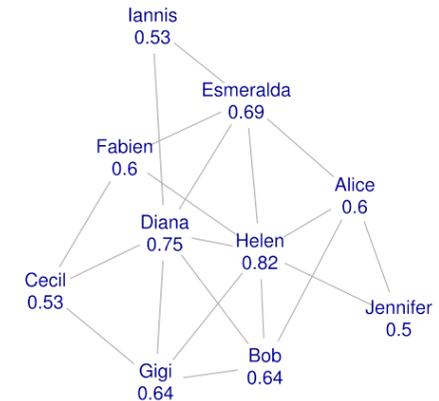
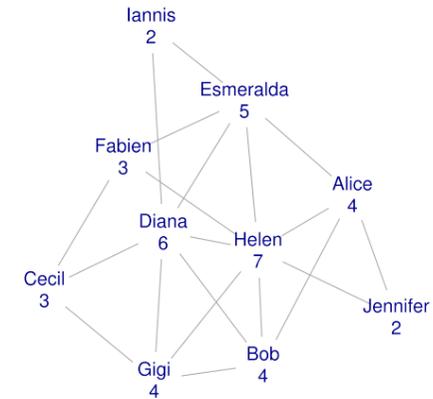
Packages: *igraph*:

- Nice functions that you cannot find anywhere else
- Uses Sparse Matrices
- Implemented in C
- Some support for bipartite networks

Rmysql, Matrix (sparse m)

Centrality

| Method | Formula | Time complexity |
|-------------|---|---|
| Degree | $C_d(v) = \frac{deg(v)}{n - 1}$ <p>Where $deg(v)$ is the number of connections that v has.</p> | $O(E)$ |
| Closeness | $C_v = \frac{ v - 1}{\sum_{i \neq v} d_{vi}}$ <p>Where d_{vi} is the distance between vertex d and i and v is the number of vertices</p> | $O(V^3)$ |
| Betweenness | $B_v = \sum_{i \neq j, i \neq v, j \neq v} g_{ivj} / g_{ij}$ <p>Where g_{ivj} is the number of shortest paths between i and j that pass through v, and g_{ij} is the number of paths between i and j that do not go through v</p> | $O(VE)$ time using Brandes' (2001) algorithm; parallelizable (Bader & Madduri, 2006). |



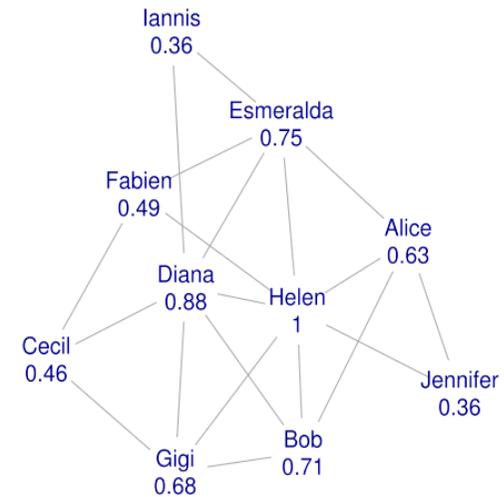
Centrality

Eigenvector

$O(V)$

$$E_v = \frac{1}{\lambda} A_{iv} E_i, \quad Ax = \lambda x$$

Where A is a matrix that represents a linear transformation, and lambda is the scaling factor (eigenvalue) in the eigenvalue equation (right)



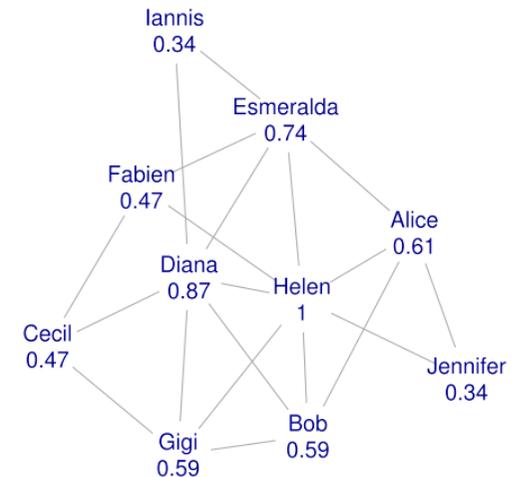
PageRank

$O(|E|\log(1/\epsilon))$, ϵ is the precision required.

$$E_v = \frac{1-d}{|V|} + d \sum_{i=1}^{|V|} A_{iv} E_i$$

The complexity is independent of the number of vertices (Bianchini, 2005, p. 100).

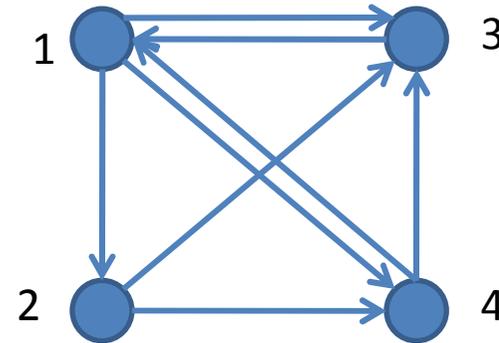
d is the damping factor (set at .85 in our experiment, as in the original Brin and Page paper)



Pagerank

- The pagerank vector is the stationary distribution of a markov chain in a link matrix
- Some assumptions to warrant convergence
- The typical value of d is .85

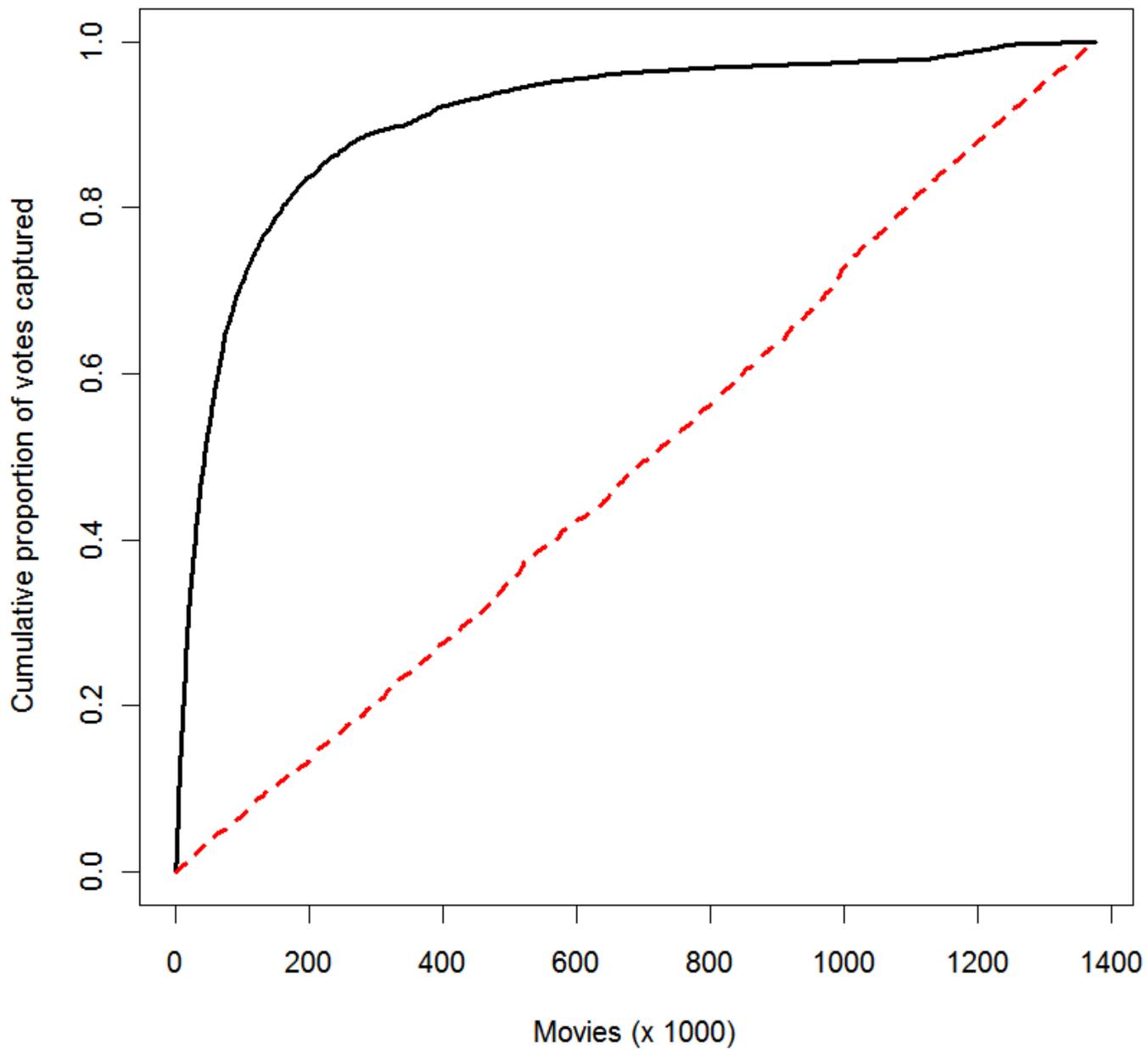
$$E_v = \frac{1-d}{|V|} + d \sum_{i=1}^{|V|} A_{iv} E_i$$



$$A = \begin{bmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{bmatrix}.$$

```
norm <- function(x) x/sum(x)
norm(eigen(0.15/nVertices + 0.85 * t(A))$vectors[,1])
```

Proportion of interest captured when subsetting
by pagerank



Top movies by pageRank in the actor->movie network

| degree | pagerank | cluster | imdbID | title | rank | votes |
|--------|-----------------------|---------|---------|--|-------|-------|
| 1298 | 0.000243688252192870 | 0 | 822609 | Around the World in Eighty Days (1956) | 40031 | 6134 |
| 313 | 0.000103540862390464 | 0 | 76352 | \Beyond Our Control\" (1968)" | 0 | 0 |
| 291 | 0.0000916690099912811 | 0 | 993780 | Gone to Earth (1950) | 7.0 | 291 |
| 285 | 0.0000890255923652847 | 0 | 915626 | Deadlands 2: Trapped (2008) | 39971 | 15 |
| 424 | 0.000083882328163772 | 0 | 1282574 | Stuck on You (2003) | 6.0 | 19709 |
| 629 | 0.0000808241101098043 | 0 | 622100 | \Shortland Street\" (1992)" | 39850 | 225 |

Problems

- Graphs have advantages over RDBMS/tables[1]. But we are used to think in tables
- There is no direct way to handle RDF in R. worth an R package?

ActiveRDF: Object-Oriented Semantic Web Programming

Eyal Oren

eyal.oren@deri.org

Renaud Delbru

renaud.delbru@deri.org

Sebastian Gerke

sebastian.gerke@deri.org

Armin Haller

armin.haller@deri.org

Stefan Decker

stefan.decker@deri.org

Digital Enterprise Research Institute
National University of Ireland, Galway
Galway, Ireland



Linked data are out there for the grabs

We need to start thinking in terms of graphs,
and slowly move away from tables

Thanks for your attention

Jose Quesada, quesada@workingcogs.com, <http://josequesada.name>

Twitter: @Quesada