Is this really Science? Evaluating research contributions

Abraham Bernstein (University of Zurich) Natasha Noy (Google Research)

tps://en.wikipedia.org/wiki/Volta_Laboratory_and_Bureau#/media/File:Photographing_Sound_in_1884._A_rare_photograph_taken_at_Volta_Laboratory_by_J._Harris_Rogerer_a_friend_of_Bell_and_Tainter_(Smithsonian_photo_44312-E)_1009.jpg

Science = Missenschaft (German)

Wissen = Knowledge Schaffen = Creating

Science = Knowledge Creation

standing on the shoulders of Giants

attributed to John of Salisbury

discovering truth by building on previous discoveries

nominalism

Ontology

Is realty objective or a product of peoples' minds?

realism

positivism

Epistemology

Do regularities exist apart from individuals?

determinism

Human Nature

What is the nature of human existence?

anti-positivism

voluntarism

Burrell and Morgan, 1979

positivism

Epistemology

Do regularities exist apart from individuals?

anti-positivism

- There is a reality!
- There is a truth
- Science is learning the rules that govern reality

Structuration Theory

- Reality is a social construct!
- Truth is contextdependent
- Science is learning the regularities within the relevant contexts

Structuration Theory



https://pixabay.com/en/videos/crosswalk-crowd-people-many-3630/

Who decides what is true and relevant? How do I convince them?



scientific Community

Nature,

Science

Chemistry Journal

Physical Chemistry Journal

Journal of Special Physical Chemistry



Truch in Science

Who decides what is true and relevant? How do I convince them?

convince... of what? What kind of Questions do we have?

All Science is either Physics or Stamp Collection Ernest Rutherford





What is a Research Question?

Research question

Determines where and what kind of research the writer will be doing Identifies the specific objectives the study or the paper will address

What is a Research Question?

Qualitative Template:

(How or what) is the ______ ("story for" for narrative research; "meaning of" the phenomenon for phenomenology; "theory that explains the process of" for grounded theory; "culture-sharing pattern" for ethnography; "issue" in the "case" for case study) of ______ (central phenomenon) for _______ (participants) at ______ (research site).

Quantitative Template:

Does ______ (name the theory) explain the relationship between (independent variable) and ______ (dependent variable), controlling for the effects of ______ (control variable)?

Creswell's (2009)



RQ and different methods...

Feasibility study

- Case study (aka Demonstrator)
- Comparative study / Benchmark
- Observational Study [a.k.a. Ethnography]
- Series Experiment
- Literature survey (incl. Meta-Analysis)Formal Model
- Simulation

The Actual Semantic Web Research Projects



http://www.deviantart.com/art/The-Good-The-Bad-and-The-Ugly-320626352

"My Semantic Web system is better than your Semantic Web system"

> What does "better" mean and how do you measure it?

Improve performance of a system



Hypothesis: "My reasoner is very fast and efficient"
 Nothing to measure here. You



Hypothesis: "My reasoner is faster than a reasoner X on one specific ontology" returnal validity:



Hypothesis: "Using X will improve the efficiency of reasoning on the class of languages Y, compared to the current state of the art"

how important is this one ontology?

Evaluation Methods (Internal validity)

Run your reasoner on large ontologies in the class Y

Run the best existing reasoner that is designed for the class of ontologies that you consider

 Run experiments to understand why it is better

 Compare the performance of your reasoner to the existing one(s)

gold standard, published benchmark



Make sure that your hypothesis is task-specific: X is better than Y for task Z (or in a context C)

- Maybe design a hierarchy of hypotheses (unfortunately, not very common in CS/AI/SemWeb)
- Make sure that your evaluation is designed to compare X and Y in the context C or for task Z
- When you report results and reach conclusions, do not over-generalize. The conclusions are valid only for these tasks/context.

What do others need to stand on your shoulders?

"Look, Ma, no hands!" or "We built a system"



Where is a <u>scientific</u> problem in building an application?

Developing an application



"I have developed a tool. It works!" Nothing to measure.



Hypothesis: "My application works perfectly for displaying one ontology and if we ask users questions about this ontology, they can use the tool effectively"

External validity:

how important is this one ontology?

No way it can fail.



 Hypothesis: "Domain experts can use our system effectively to accomplish a task X (e.g., map between large ontologies)"



 Hypothesis: "Domain experts can use our system more effectively than another system
 Z to accomplish a task X (e.g., map between large ontologies)"

Evaluation Methods

Experiment
Usability study
Successful completion of tasks
Case study
Comparative study or benchmark



- Just developing a system is not a research contribution in itself.
- Make sure that your hypothesis is task-specific: X is good for task Z, or in context C, or for users of type T
- Maybe design a hierarch.
 (unfortunate We have seen this before!
- Make sure that your evaluation is designed to compare X and Y task Z, or in context C, or for users of type T
- When you report results and reach conclusions, do not overgeneralize. The conclusions are valid only for these tasks/context/ user types.

We will put everything in RDF and the world will be a better place

Solution in search of a problem: who should care if you succeed?

Convert unstructured data into RDF or Linked Data



• "I will convert a corpus of abstracts into RDF"



Hypothesis: "My conversion process produces better linked data than conversion process X"
Let who cares?



Hypothesis: "Using extracted linked data will improve search performance on the corpus compared to existing methods"



 Hypothesis: "Using extracted linked data will enable advanced querying that was not possible before"
 Make sure there is somebody who



 actually wanted on the corpus queries on the corpus queries on the corpus only one): "My method for extracting structured data has better accuracy/ coverage/precision/recall/etc. than the state of the art." "Using our hammer for every nail" does the use of semantic web technology actually improve anything?



Novel Solution to an Old Problem



"We will use Semantic Web technology to make movie recommendations"

Key difference: Our goal is to improve the efficiency of a task that someone cares about. Not to use SW technology per se



Hypothesis: "We will improve the efficiency of social-network monitoring by using SW technology/improve the quality of recommendations"

Evaluation Methods

Compare the accuracy of recommendations with and without the linked data component

Compare the accuracy of your system to an existing non-Semantic Web system



If the LD component is completely integral to your system and you cannot take it out, you will need to compare to another system

- You may need to compare to the state of the art to convince non-SemWebbies that your method has any value
- Make sure the metrics, the users, and the datasets are comparable
- Think how others can re-use your results
 - What if your testbed is the LOD Cloud?

What (human) languages are ontologies published in?

> Stamp collection for the sake of stamp collection

http://www.flickr.com/photos/ rachelfordjames/2833420148

Improve performance of a system



o "We will create a set of features that ontologies have and will describe ontologies from our catalog according to these features" This is not a hypothesis; this



o "We will create a set of features that ontologies have and will describe ontologies from representative ontology repositories according to these features."

Why would this information be useful? What will drive you selection of features? Imagination?

is a research plan

Improve performance of a system



Hypothesis: "Only a small number of OWL constructs are used in the publicly available ontologies."

Why should anyone care?

Example: if a particular set of features are almost never used, so it can be ok that your reasoner does not support it.

Evaluation Methods

Collect a representative corpus / data
<u>Representative</u> is the operative word here
Analyze the terms used and determine which ones are not used much

Reproducibility:

Extremely difficult in hypotheses generating studies!



Develop a formal model/ workflow/etc for a task



 "We will develop a workflow for creating linked data and annotating it with ontologies"
 Nothing to measure here. You



We will study a number of existing workflows and will create a more general one."



 Hypothesis: "It is possible to build a formal workflow for collaborative creation of linked data (similar: It is possible to develop a formal representation for X)"

Not great: You cannot falsify this



Hypothesis: "My workflow model is generic enough to represent a meaningful number of diverse published workflows".



Auxiliary hypothesis: "My system provides sound and complete reasoning."



 Auxiliary hypothesis: "My formalism elements are symmetric, reflexive, transitive."

Evaluation Methods

Prove a theorem!

Find a <u>representative</u> set of workflows/ problems/etc and represent in your model

Re-cap: Types of problems

"My Semantic Web system is better than your Semantic Web system" "Look, Ma, no hands!" or "We built a system"

How can we use Linked data to solve this problem?



We will put everything in RDF and the world will be a better place

Stamp collection for the sake of stamp collecting



What have we learned?

Make sure

- o you have a good/appropriate research questions
- you operationalized your research questions with (falsifiable) hypotheses
- your evaluation plan is designed to test your hypothesis.
- "Who cares?" and "So what?"
- Always think about the limitations / threats to validity!

Main Take-Aways





- What is your RQ?
- @ "Who cares?" and "So what?"
- Always think about the limitations /
 threats to validity!

Where to go from here?

tion

Charles M. Judd

Eliot R. Smith

Louise H. Kidder



No. IFI-2014.02

1

TECHNICAL REPORT

Abraham Bernstein Natasha Noy

Is This Really Science? The Semantic Webber's Guide to Evaluating Research Contributions

April 2014

University of Zurich Department of Informatics (IFI) Binzmühlestrasse 14, CH-8050 Zürich, Switzerland

