Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

Constructing a Knowledge Base on Aging An Automated Approach

Mark Farrell

Undergraduate Student, University of Waterloo

Center for Research and Education on Aging Lawrence Berkeley National Laboratory University of California, Berkeley

September 4th, 2014

Outline

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

1 Automatically Constructing Knowledge Bases

2 Extracting Facts in a Structured Forma

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussior ■ CREA is constructing a knowledge base to study and understand the human aging process.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- CREA is constructing a knowledge base to study and understand the human aging process.
- New discoveries are published quickly and in large volume.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- CREA is constructing a knowledge base to study and understand the human aging process.
- New discoveries are published quickly and in large volume.
- It is infeasible to construct the knowledge base by hand.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- CREA is constructing a knowledge base to study and understand the human aging process.
- New discoveries are published quickly and in large volume.
- It is infeasible to construct the knowledge base by hand.
- Working on software to construct the knowledge base automatically.

How to Automatically Construct the Knowledge Base

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured

Results &

■ Routinely search for keywords related to aging, dowloading text articles from sources like PubMed and WebMD.

How to Automatically Construct the Knowledge Base

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting
Facts in a
Structured
Format

- Routinely search for keywords related to aging, dowloading text articles from sources like PubMed and WebMD.
- Build a spam filter to get rid of non-scientific sentences.

How to Automatically Construct the Knowledge Base

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Routinely search for keywords related to aging, dowloading text articles from sources like PubMed and WebMD.
- Build a spam filter to get rid of non-scientific sentences.
- Extract scientific facts from the sentences and save them in a structured format.

How to Automatically Construct the Knowledge Base

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Routinely search for keywords related to aging, dowloading text articles from sources like PubMed and WebMD.
- Build a spam filter to get rid of non-scientific sentences.
- Extract scientific facts from the sentences and save them in a structured format.
- Provide a graphical interface that allows users to search and otherwise explore the knowledge base.

Summary of Progress

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

 Devised and implemented the method for finding simple facts in sentences, extracting them in a structured format.

Summary of Progress

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured

- Devised and implemented the method for finding simple facts in sentences, extracting them in a structured format.
- Began work on a web viewer for the knowledge base.

Outline

Constructing a Knowledge Base on Aging

Mark Farrell

Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

1 Automatically Constructing Knowledge Bases

2 Extracting Facts in a Structured Format

Tokenization

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

■ Input a text document and read it, one sentence at a time.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

Example: Tokenization

scala> tokens("The man walks. The dog eats.") res0: List[String] = List(The man walks., The dog eats.)

Parsing

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

■ For each sentence, generate a constituent tree that describes its phrase structure.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

Example: Parsing

```
scala> parse("The man walks the dog.")
res0: Tree[String] = (ROOT
  (S
     (@S
          (NP (DT The) (NN man))
          (VP (VBZ walks)
                (NP (DT the) (NN dog))))
          (. .)))
```

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion ■ The University of Pennsylvania Treebank Project:

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- The University of Pennsylvania Treebank Project:
 - Defines notation for constituent trees.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- The University of Pennsylvania Treebank Project:
 - Defines notation for constituent trees.
 - Parses sentences from the Wall Street Journal by hand.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- The University of Pennsylvania Treebank Project:
 - Defines notation for constituent trees.
 - Parses sentences from the Wall Street Journal by hand.
- The Berkeley Parser is software that guesses how to parse a sentence from the notation and examples specified by the Penn Treebank.

Compilation

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

■ Extract facts from each constituent tree.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

Example: Compilation

scala> compile("The man walks the dog.").shows res0: String = [<compound:walk(<atom:man>, <atom:dog>)>]

Compilation Method

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion Pattern match on the constituent trees. Define patterns for:

1 Extracting nouns from noun phrases (NP).

Compilation Method

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

Pattern match on the constituent trees. Define patterns for:

- 1 Extracting nouns from noun phrases (NP).
- 2 Extracting predicates and nouns from verb phrases (VP).

Compilation Method

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion Pattern match on the constituent trees. Define patterns for:

- 1 Extracting nouns from noun phrases (NP).
- **2** Extracting predicates and nouns from verb phrases (VP).
- 3 Extracting facts from complete clauses (S), making logical assertions with nouns and predicates.

Outline

Constructing a Knowledge Base on Aging

Mark Farrell

Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

1 Automatically Constructing Knowledge Bases

2 Extracting Facts in a Structured Format

Software Demonstration

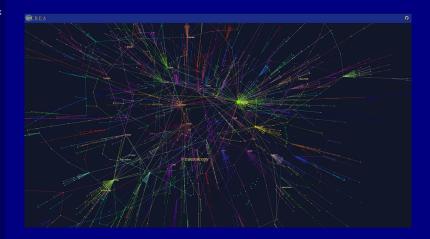
A preview of CREA's knowledge base, compiled from PubMed abstracts.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format



Performance Parallelization

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion ■ It is possible to extract facts from many sentences at the same time.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results & Discussion

■ Filter spam sentences from documents.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Filter spam sentences from documents.
- The accuracy of the parser could be optimized:

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Filter spam sentences from documents.
- The accuracy of the parser could be optimized:
 - Should be trained to identify more nouns from the biomedical domain.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Filter spam sentences from documents.
- The accuracy of the parser could be optimized:
 - Should be trained to identify more nouns from the biomedical domain.
- Define more patterns for extracting facts:

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Filter spam sentences from documents.
- The accuracy of the parser could be optimized:
 - Should be trained to identify more nouns from the biomedical domain.
- Define more patterns for extracting facts:
 - The software succeeds around 50% of the time.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

■ Support negated clauses and conditional logic.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Support negated clauses and conditional logic.
- Facts can contradict each other:

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Support negated clauses and conditional logic.
- Facts can contradict each other:
 - Store the probability that is true as the weight of its edge on the knowledge base's graph.

Constructing a Knowledge Base on Aging

Mark Farrell

Automaticall Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Support negated clauses and conditional logic.
- Facts can contradict each other:
 - Store the probability that is true as the weight of its edge on the knowledge base's graph.
- Scale and launch the software service.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

Results &

Demonstrated a method for automatically constructing CREA's knowledge base on aging.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Demonstrated a method for automatically constructing CREA's knowledge base on aging.
- Showed how to extract facts from English text in the knowledge base's structured format.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Demonstrated a method for automatically constructing CREA's knowledge base on aging.
- Showed how to extract facts from English text in the knowledge base's structured format.
- Discussed the need to improve software accuracy by lensing in on the biomedical domain.

Constructing a Knowledge Base on Aging

Mark Farrell

Automatically Constructing Knowledge Bases

Extracting Facts in a Structured Format

- Demonstrated a method for automatically constructing CREA's knowledge base on aging.
- Showed how to extract facts from English text in the knowledge base's structured format.
- Discussed the need to improve software accuracy by lensing in on the biomedical domain.
- Suggested how the software implementation can be scaled for production usage.