# Wrapper Induction

Fabian M. Suchanek

# Semantic IE

Reasoning

You are here

Fact Extraction

Instance Extraction

singer

Entity Disambiguation

singer Elvis

Entity Recognition

Source Selection and Preparation
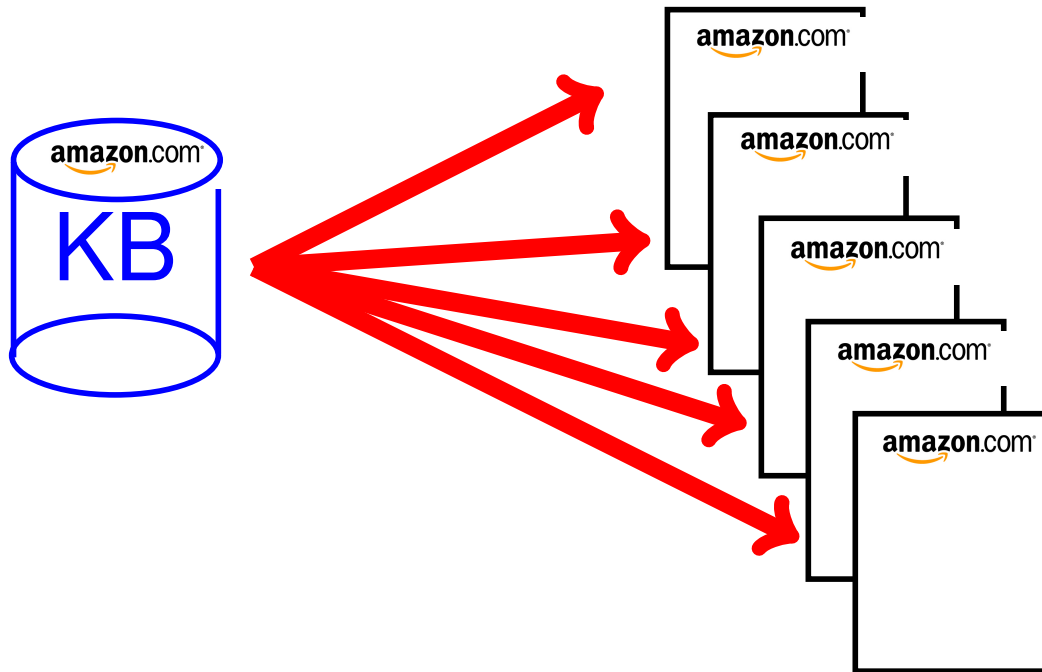
2

# Generated Web pages
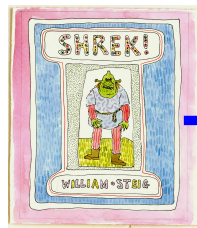
Web page generation is the process of producing several similar Web pages from a KB.

# Example: Generated Web pages



amazon.com

price → 10 USD

author → WilliamSteig

amazon.com

"Shrek"

by W. Steig

only 10 USD!

Buy it!

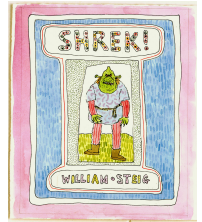# IE aims to reconstruct the KB



price → 10 USD

writer → W.Steig

amazon.com

"Shrek"

by W. Steig

only 10 USD!

Buy it!

5

# Def: Wrapper

A wrapper for a set of pages generated

from the same KB is a function that

extracts strings from such a page.

(Technically, it is the inverse function of the function that generated the page.
The strings still have to be disambiguated and put in relation to yield facts.
Different applications have different more specific definitions of the "strings".)
Kushmerick: Wrapper Induction



"Shrek…",

"90 min",

"7.9"

# Information is always in same place



If we understand this...

then we understand this.

# Positions in Web pages

How do we identify

a position in a Web page?

Let's consider a simple method.

# Def: XPath

XPath is a formal language for selecting

nodes in an XML document. Wikipedia/XPath

/ identifies the root node

K/T[i] identifies the i-th child with

tag T of the node identified by K

K/T is K/T[1] if K has one T child

(XPath has many more expressions, allowing the union of node sets, the selection by attributes, etc.)

# Example: XPath

```
<html>
  <body>
    <h1>Aloha from Hawaii</h1>
    <p>This is a really great movie</p>
    <p>Stars:<i>Elvis Presley</i></p>
  </body>
</html>
```

Try it out

/html/body/p[2]/i

(Technically, this expression identifies the <i> node, not the text inside.)

# Task: XPath

Write XPath expressions that identify
nodes whose text is "Shrek", "W. Steig",
and "84 min".

```
<html>
  <body>
    <b>Shrek</b>
    <ul>
      <li>Creator: <b>W. Steig</b>
      <li>Duration: <i>84 min</i>
    </ul>
```

11

# Def: Wrapper induction

Wrapper induction is the process of generating a wrapper from a set of Web pages with strings to be extracted.



+

"Shrek", "7.9"

=

/html/body/h1
/html/body/p[2]/i

Web page

+

Strings to

be extracted

=

Wrapper

# Wrapper induction

Wrapper Induction requires as input

Web pages with strings to be extracted.

These can come, e.g.,

- from a KB

$$hasTitle(ShrekMovie, "Shrek")$$

- from manual extraction

 + 👤 = "Shrek"

- from manual annotation in a GUI

 + 👤 =

# Def: Wrapper Application

Wrapper application is the process of extracting its strings from a Web page.



**Web page**

**+**

/html/body/h1
/html/body/p[2]/i

**+**

**Wrapper**

**=**

**=**

"Elvis", "11"

**Strings**

# From Strings to Facts

The extracted strings have to be disambiguated and put into relations in order to yield facts about entities.



"Aloha from Hawaii"

"Elvis Presley"

"11"

hasTitle(e42, "Aloha from Hawaii")

hasActor(e42, ElvisPresley)

hasRating(e42, "11.0")

advanced>

# Detail pages

A detail page contains information about one entity, the "page entity".

# List pages

List pages contain information on several entities.

# Data may exhibit structure

Dronkeys:

&lt;ul&gt;

&lt;li&gt;Eclair: female

&lt;li&gt;Bananas: flexible

&lt;/ul&gt;

Shrek's kids:

&lt;ul&gt;

&lt;li&gt;Farkle: male

&lt;li&gt;Fergus: male



one question.
HOW? vi.sualize.us

# (Web page) Types

A type is a name plus one of

- a basic type

- a set of types

- a tuple of types

family: tuple (

  name: string

  children: set (

    child: tuple (name: string,

          gender: string)))

# ROADRUNNER: Learn types

Page 1:

<ul>

<li>Peanut

</ul>

Page 2:

<ul>

<li>Charles

</ul>

Wrapper:

<ul>

<li>[FIELD]

</ul>

# ROADRUNNER: Learn types

Page 1:

```
<ul>
<li>Peanut
</ul>
```

Page 2:

```
<ul>
<li>Charles
<li>Anne
</ul>
```

Wrapper:
```
<ul>
(<li>[FIELD])+
</ul>
```

(Set type)

advanced>

ROADRUNNER

21

# Data not separated by tags

\<html\>

\<body\>

 The Dronkeys\<br\>

 \<p\>

Eclair: female\<br\>

Bananas: male\<br\>

Peanut: unknown\<br\>

 \<hr\>

# WIEN: Wrappers in HTLR language

For the WIEN system, a wrapper is of the form

$$\langle head, tail, left, right, left, right, ...\rangle$$

e.g.

$$\langle\ \langle p\rangle, \langle hr\rangle, n, :, , \langle br\rangle\rangle$$

# Applying a HTLR wrapper

$$< <p>,<hr>,n , : , , <br>>$$

head    tail    left    right    left    right

<html>

<body>

 The Dronkeys<br>

 <p>↓

 Eclair: female<br>

Bananas: male<br>

Peanut: unknown<br>

<hr>

1. scroll

   to head

24

# Applying a HTLR wrapper

< <p>,<hr>,n , : , , <br>>
head    tail    left    right    left    right

<html>

<body>

  The Dronkeys<br>

  <p>

Eclair: female<br>

Bananas: male<br>

Peanut: unknown<br>

<hr>

2. scroll
to left

# Applying a HTLR wrapper

$<$ $<p>$, $<hr>$, n , : , , $<br>$ $>$

head     tail     left     right     left     right

$<html>$

$<body>$

The Dronkeys$<br>$

$<p>$

Eclair: female$<br>$

Bananas: male$<br>$

Peanut: unknown$<br>$

$<hr>$

3. extract

until

right

$<$"Eclair",

26

# Applying a HTLR wrapper

$$< <p>,<hr>,n , : , , <br>>$$

head     tail     left     right     left     right

<html>

<body>

  The Dronkeys<br>

  <p>

  Eclair: female<br>

  Bananas: male<br>

  Peanut: unknown<br>

  <hr>

4. scroll

to left

<"Eclair",

# Applying a HTLR wrapper

$$< <p>, <hr>, n, :, , <br>>$$

head     tail     left     right     left     right

<html>

<body>

  The Dronkeys<br>

  <p>

Eclair: female<br>

Bananas: male<br>

Peanut: unknown<br>

<hr>

5. extract

until

right

<"Eclair",

"female"

28

# Applying a HTLR wrapper

$$< <p>,<hr>,n , : , , <br>>$$

head     tail     left     right     left     right

$<$html$>$

$<$body$>$

  The Dronkeys$<$br$>$

  $<$p$>$

Eclair: female$<$br$>$

Bananas: male$<$br$>$

Peanut: unknown$<$br$>$

$<$hr$>$

6. close

  the tuple

$<$"Eclair",

  "female"$>$

# Applying a HTLR wrapper

$< \ <p>,<hr>,n \ , : \ , \ , \ <br>>$

head     tail     left     right     left     right

<html>

<body>

     **7. repeat**

 The Dronkeys<br>

         **until tail**

 <p>

Eclair: female<br>

Bananas: male<br>     <"Eclair", "male">

    <"Bananas","male">

Peanut: unknown<br>     <"Peanut","unknown">

                             advanced>

<hr>

# Delimiters may differ

&lt;html&gt;

&lt;body&gt;

The Dronkeys&lt;br&gt;

&lt;p&gt;

Eclair: female&lt;br&gt;

Bananas (male)&lt;br&gt;

Peanut is hybrid&lt;br&gt;

&lt;hr&gt;

# STALKER: Disjunctions in Wrappers

Start of gender:

- punctuation

- or "is"

Eclair: female$<$br$>$

Bananas (male)$<$br$>$

Peanut is hybrid$<$br$>$

# Overview

- WordNet
- Fact Extraction
  - Fact Extraction from Wikipedia
    - Infoboxes
    - Categories
    - Checks
  - Fact Extraction from generated pages
    - Simple Wrappers
    - Advanced Wrappers
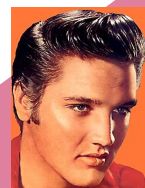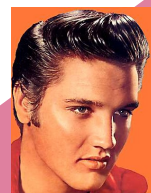
# Semantic IE

Reasoning



You are here

Fact Extraction

Is-A Extraction

singer

Entity Disambiguation

singer Elvis

Entity Recognition

Source Selection and Preparation

# References

Kushmerick: Wrapper Induction

ROADRUNNER

Web data mining class

Muslea: STALKER