# Non IT congurable adaptive data mining solution used in transforming raw data to structured data

**Bachelor's Thesis in Artificial Intelligence**
**Radboud University Nijmegen**

**Franc Grootjen**
**RU**

**Alessandro Paula**
**Hyperleap**

Mart Lubbers
s4109053

July 24, 2014

# Contents

**Abstract**

## Abstract

# 1 Introduction

## 1.1 Introduction

Within the entertainment business there is no consistent style of informing people about the events. Different venues display their, often incomplete, information in entirely different ways. Because of this, converting raw information from venues to structured consistent data is a challenging and, relevant problem.
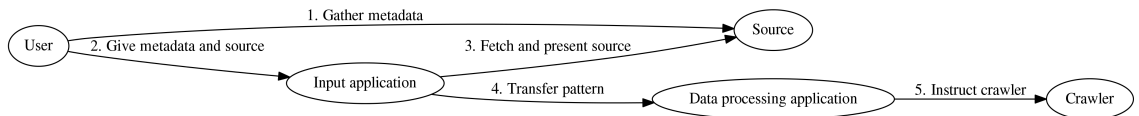
## 1.2 HyperLeap

Hyperleap is a small company that is specialized in infotainment (information + entertainment) and administrates several websites which bundle information about entertainment in an ordered way and as complete as possible. Right now, most of the input data is added to the database by by hand which is very labor intensive. Therefore Hyperleap is looking for a smart solution to automate a part of the data injection in the database, the crux however is that the system must not be too complicated from the outside and be useable for a non IT professional(NIP).

## 1.3 Research question and practical goals

This brings up the main research question: *How can we make an adaptive, autonomous and programmable data mining program that can be set up by a NIP which is able to transform raw data into structured data.*

In practice the goal and aim of the project is to create an application that can, with NIP input, give computer parseable patterns which a separate crawler can periodically crawl. The NIP has to be able to enter the information about the data source in a user friendly interface which sends the information together with the data source to the data processing application. The dataprocessing application then in turn processes the data into a extraction pattern which is sent to the crawler. The crawler can visit sources specified by the NIP accompanied by the extraction pattern created by the dataprocessing application. This workflow is described in graph 1.1.

Figure 1.1: Workflow within the applications



In this way the NIP can train the crawler to periodically crawl different data sources without too much technical knowledge. The main goal of this project is to extract the underlying structure rather then to extract the substructures. The project is in principle a continuation of a past project done by Wouter Roelofs[**?**] which was also supervised by Franc Grootjen and Alessandro Paula, however it was neven taken out of the experimental phase. The techniques described by Roelofs et al. are more focussed on extracting data from substructures so it can be an addition to the current project.

As a very important sidenote, the crawler needs to notify the administrators if a source has become problematic to crawl, in this way the NIP can very easily retrain the application to fit the latest structural patterns.

## 1.4  Scientific relevance

Currently the techniques for conversion from non structured data to structured data are static and mainly only usable by IT specialists. There is a great need of data mining in non structured data because the data within companies and on the internet is piling up and are usually left to catch dust.

# 2 Methods

## 2.1 Input application

The user input all goes through the familiar interface of the user's preferred web browser. By visiting the crawler's train website the user can specify the metadata of the source it wants to be periodically crawled through simple web forms as seen in figure 2.1

Figure 2.1: Webforms for source metadata

Venue:

Frequency:

Default location name:

Adress:

Website:

## 2.2 Data processing application

### 2.2.1 Directed acyclic graphs and finiti automata

Directed acyclic graphs(DAG) and finite state automaton(FSA) have a lot in common concerning pattern recognition and information extraction. By feeding words into an algorithm a DAG can be generated so that it matches certain patters present in the given words. Figure 2.2 for example shows a FSA that matches on the words *ab* and *ac*.

Figure 2.2: Example DAG/FSA



With this FSA we can test if a word fits to the constraints it the FSA describes. And with a little adaptation we can extract dynamic information from semi-structured data.

### 2.2.2 Back to DAG's and FSA's

Nodes in this data structure can be single letters but also bigger constructions. The example in Figure 2.3 describes different separator pattern for event data with its three component: what, when, where. In this example the nodes with the labels *what, when, where* can also be complete subgrahps. In this way data on a larger level can be using the NIP markings and data within the categories can be processed autonomously.

Figure 2.3: Example event data



### 2.2.3 Algorithm

## 2.3 Crawler application

# 3 Results

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim

interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetuer.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

# 4 Discussion

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim

interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetuer.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

# 5 Appendices

## 5.1 Input application

<div align="center">Python front/back-end</div>

```
1    #!/bin/env python
2    # −∗− coding: utf−8 −∗−
3
4    from mod_python import apache, util
5    import feedparser
6    import re
7    import urllib
8    import subprocess
9    import os
10
11   def req_pre_pos(req):
12       req.log_error('handler')
13       req.content_type = 'text/html'
14       req.send_http_header()
15       args = util.FieldStorage(req)
16       req.write(""" \
17   <html>
18   ____<head>
19   _____<title>VER:_0.01_−_HyperFrontend_RSS_feed_POSTREQUEST</title>
20   ____</head>
21   ____<body>
22   _____Thanks_submitting:_<br_/>
23   _____<a_href="index.html">Enter_new_rss_feed</a>
24   _____<pre>
25   {}
26   _____</pre>
27   ____</body>
28   </html>
29   """.format(args))
30       os.chdir('/var/www/py/files')
31       with open('raw_out.txt', 'w') as f:
32           f.write(str(args))
33
34
35   def req_pre(req):
36       req.log_error('handler')
37       req.content_type = 'text/html'
38       req.send_http_header()
39       req.write(""" \
40   <html>
41   ____<head>
42   _____<title>HyperFrontend_RSS_feed_input</title>
43   _____<script_src="contextmenu_o.js"></script>
44   ____</head>
45   ____<body>
46
47   ____<table>
```

```
48          <tr><td>Venue: </td><td>
49              <input type="text" name="venue" class="indat"></td></tr>
50          <tr><td>Frequency: </td><td>
51              <input type="text" name="freq" class="indat"></td></tr>
52          <tr><td>Default location name: </td><td>
53              <input type="text" name="dloc" class="indat"></td></tr>
54          <tr><td>Adress: </td><td>
55              <input type="text" name="adress" class="indat"></td></tr>
56          <tr><td>Website: </td><td>
57              <input type="text" name="website" class="indat"></td></tr>
58      </table>
59
60      Selecteer iets en klik de link<br/>
61      <button style="color:blue" onclick="javascript:f_waar()">Waar</button>
62      <button style="color:green" onclick="javascript:f_wat()">Wat</button>
63      <button style="color:red" onclick="javascript:f_wanneer()">Wanneer</button>
64      <br/>
65
66  <div style="position:absolute; left :12px;width:500px;"></div>
67  <script language="javascript" type="text/javascript">
68      var content='<b>Categorize</b><br/>';
69      content+='<a href="#" onclick="javascript:f_waar()">Waar</a><br/>';
70      content+='<a href="#" onclick="javascript:f_wat()">Wat</a>';
71      content+='<a href="#" onclick="javascript:f_wanneer()>Wanneer</a><br/>';
72      content+=' Test 123';
73      init (content,120);
74  </script>
75  """)
76
77
78  def req_post(req):
79      req.write("""\
80              <button onclick="javascript:post_all()" method="post" target="_bla\
81  nk">Submit</button>
82      </body>
83  </html>
84  """)
85
86
87  def feed2html(req, url, name):
88      url = urllib.unquote(url)
89      url = url if re.match('https?://', url) else 'http://{}'.format(url)
90      req.write(
91          '\tLoading "{}" as <p id="rssname">"{}"</p><br/>\n'.format(url, name))
92      feed = feedparser.parse(url)
93      channel = feed.feed
94      req.write('\t<table id="content-table" border="1" id="htab">\n')
95      req.write('\t\t<tr><th>Title</th><th>Summary</th></tr>\n')
96      for i in feed.entries:
97          req.write('\t\t<tr><td>{}</td><td>{}</td></tr>\n'.
98                      format(i['title'].encode('ascii', 'xmlcharrefreplace'),
99                              i['summary'].encode('ascii', 'xmlcharrefreplace')))
100     req.write('\t</table>\n<br/>')
101
102
103 def handler(req):
104     if req.method == "POST":
105         req_pre_pos(req)
106     else:
107         req_pre(req)
108         args = util.FieldStorage(req)
109         if 'url' not in args and 'name' not in args:
```

```
110             req.write('Something_went_wrong,_empty_fields?<br_/>')
111             req.write('<a_href="index.html">back</a>')
112         else:
113             feed2html(req, args['url'], args['name'])
114         req_post(req)
115     return apache.OK
```

## HTML landing page

```
1  <html>
2      <head>
3      </head>
4      <body>
5          <form method="get" action="./hyper.py">
6              <table>
7                  <tr><td><p>RSS URL:   </td><td><input type="text" name="url"
8                          value="localhost/py/paradiso.rss.xml"></td></tr>
9                  <tr><td>RSS Name: </td><td><input type="text" name="name"></td></
                        tr>
10                 <tr><td><input type="submit" value="Submit"></p>
11             </table>
12         </form>
13     </body>
14  </html>
```

## Javascript frontend

```
1  var selection;
2  var mouse_x = 0;
3  var mouse_y = 0;
4  var mouse_left = false;
5  var mouse_right = false;
6  if (document.addEventListener != undefined) document.addEventListener('mousemove',
        mouseMove, true);
7  else if (document.layers) document.captureEvents(Event.MOUSEDOWN | Event.MOUSEMOVE
        | Event.MOUSEUP);
8  document.onmousemove = mouseMove;
9  document.oncontextmenu = RightMouseDown;
10 document.onmousedown = mouseDown;
11 document.onmouseup = mouseUp;
12
13 function mouseMove(a) {
14     mouse_x = document.all ? event.clientX + document.body.scrollLeft : document.
            layers ? a.x + window.pageXOffset : a.clientX + window.pageXOffset;
15     mouse_y = document.all ? event.clientY + document.body.scrollTop : document.
            layers ? a.y + window.pageYOffset : a.clientY + window.pageYOffset
16 }
17
18 function RightMouseDown() {
19     mouse_right = true;
20     return false
21 }
22
23 function init(a, w) {
24   console.log(a)
25     var b = document.createElement("DIV");
26     b.id = "contextmenu";
27     if (!w) var w = 120;
28     b.style.width = w + "px";
29     var c = '<div_style="position:relative;left:5px;top:-4px;">';
30     c += a;
```

17

```
31        c += '</div>';
32        b.innerHTML = c;
33        b.style.position = "absolute";
34        b.style.left = "0px";
35        b.style.top = "0px";
36        b.style.visibility = "hidden";
37        b.style.overflow = "hidden";
38        b.style.padding = "4px";
39        b.style.backgroundColor = "#ffffff";
40        b.style.border = "1px_solid_#6a6868";
41        document.body.appendChild(b);
42        delete b
43   }
44
45   function mouseUp(e) {
46        var curselection = window.getSelection().getRangeAt(0);
47        if (curselection.endOffset - curselection.startOffset > 0)
48     selection = curselection;
49        console.log(selection)
50        if (e.which == 1) document.getElementById("contextmenu").style.visibility = "
             hidden";
51        else if (e.which == 3) mouse_right = false
52   }
53
54   function mouseDown(e) {
55        if (e.which == 3) {
56             mouse_right = true;
57             document.getElementById("contextmenu").style.left = mouse_x + "px";
58             document.getElementById("contextmenu").style.top = mouse_y + "px";
59             document.getElementById("contextmenu").style.visibility = "visible"
60        }
61   }
62
63
64   function stylizeHighlightedString(range, col)
65   {
66        var selectionContents = range.extractContents();
67        var span = document.createElement("span");
68        span.appendChild(selectionContents);
69        span.setAttribute("class","uiWebviewHighlight");
70        span.style.backgroundColor = col;
71        span.style.color = "white";
72        range.insertNode(span);
73   }
74
75   function f_wanneer() {
76        stylizeHighlightedString(selection, "red")
77   }
78
79   function f_wat() {
80        stylizeHighlightedString(selection, "green")
81   }
82
83   function f_waar() {
84        stylizeHighlightedString(selection, "blue")
85   }
86
87   function post_all() {
88        var xmlhttp = new XMLHttpRequest();
89        xmlhttp.onreadystatechange=function()
90        {
91             if (xmlhttp.readyState==4)
```

```
92              {
93                   document.write(xmlhttp.responseText);
94              }
95          }
96          var params = "content="+encodeURIComponent(document.getElementById("content-
                table").innerHTML);
97          params += "&name="+encodeURIComponent(document.getElementById("rssname").
                innerHTML);
98          var indatarray = document.getElementsByClassName('indat')
99          for (var i = 0; i<indatarray.length; i++) {
100             params += "&" + indatarray[i].name + "=" + indatarray[i].value;
101         }
102         xmlhttp.open("POST", "hyper.py", true);
103         xmlhttp.setRequestHeader("Content-type", "application/x-www-form-urlencoded");
104         xmlhttp.setRequestHeader("Content-length", params.length);
105         xmlhttp.send(params);
106     }
```