# FOREST: Focused Object Retrieval by Exploiting Significant Tag paths

Marilena Oita    Pierre Senellart

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views...

- ■ **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)

- ■ Information Extraction? extracting "values" out of a common structure (e.g., wrapper induction)

- ■ Information Retrieval? having some keywords, extract relevant data (not pages, but well-defined objects → the aim of FOREST)

TELECOM
ParisTech

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views...

- **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)
- Information Extraction? extracting "values" out of a common structure (e.g., wrapper induction)
- Information Retrieval? having some keywords, extract relevant data (not pages, but well-defined objects $\rightarrow$ the aim of FOREST)

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views...

- **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)
- **Information Extraction?** extracting "values" out of a common structure (e.g., wrapper induction)
- Information Retrieval? having some keywords, extract relevant data (not pages, but well-defined objects $\rightarrow$ the aim of FOREST)

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views. . .

- ■ **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)
- ■ **Information Extraction?** extracting "values" out of a common structure (e.g., wrapper induction)
- ■ Information Retrieval? having some keywords, extract relevant data (not pages, but well-defined objects → the aim of FOREST)

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views. . .

- **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)

- **Information Extraction?** extracting "values" out of a common structure (e.g., wrapper induction)

- **Information Retrieval?** having some keywords, extract relevant data (not pages, but well-defined objects $\rightarrow$ the aim of FOREST)

TELECOM
ParisTech

Goal: extracting interesting content from Web pages
$\approx$ eliminating boilerplate

The problem has been viewed from different point of views...

- **Text Extraction?** extracting text, yes, but not not any kind of text (e.g., BOILERPIPE)
- **Information Extraction?** extracting "values" out of a common structure (e.g., wrapper induction)
- **Information Retrieval?** having some keywords, extract relevant data (not pages, but well-defined objects → the aim of FOREST)

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

TELECOM
ParisTech

unsupervised wrapper induction: given a set of objects, infer a wrapper procedure (by grammar rules/ XPath) for extracting the data values of these objects

Typical: deep Web response pages → records (e.g., Amazon books)
Here: blogs, news, social media → objects (e.g., posts, events, tweets)

Marilena Oita   Pierre Senellart

A Day of the Dead offering in the Nunkini cemetery, in Campeche, Mexico. (Jeffrey Becom/LPI)

**Related**



**Worldwide weird: Halloween to the extreme**
Taking normally odd rituals one step further



**Oktoberfest: Then and now**
It has not changed that much over time



**On the last night of the autumn harvest, the night changes from the sunny warmth of summer to the cold dark of winter, the land from fertile to barren. The ancient Celts believed this transition gave supernatural forces a chance to break through into the world of the living, and their evil mischief to flourish.**

They came to celebrate the night leading into winter as Samhain (meaning "summer's end"), the festival widely considered to be the precursor of Halloween. On Samhain night, the Celts believed, the spirits of people who had died in the past year would walk among the living, so. villagers put out food and sweets to pacify these spirits – a ritual that may have preceded trick-or-treating. (There is no hard evidence, however, that Samhain was indeed a festival of the dead, points out historian Nicholas Rogers, in his book Halloween: From Pagan Ritual to Party Night.)

Although Halloween has pagan origins, its name is derived from the Christian holiday "All Hallows Eve", or the evening before All Saints' Day (1 November). The holiday itself was adapted by Christians who hoped to stamp out paganism, and over the years, some of the darker aspects of Halloween have been replaced by more light-hearted, family-friendly festivities. But Halloween's ties with the scary and supernatural still hold strong today, in celebrations all over the world.

**Ireland**
In Ireland, arguably the holiday's birthplace, Halloween is still greeted

"Halloween, past and present"

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

- Model a Web page as its DOM (Document Object Level) tree
- Distinguish significant content nodes
    - textual DOM leaf nodes
    - at least one keyword
- Keywords automatically acquired:
    - *Tf-Idf* analysis
        - IN: set of sample Web pages
        - OUT: top-$k$ tf-idf weighted terms
    - basic text preprocessing on feed item metadata to identify top-$k$ feed keywords

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

- Model a Web page as its DOM (Document Object Level) tree
- Distinguish significant content nodes
  - textual DOM leaf nodes
  - at least one keyword
- Keywords automatically acquired:
  - *Tf-Idf* analysis
    - IN: set of sample Web pages
    - OUT: top-$k$ tf-idf weighted terms
  - basic text preprocessing on feed item metadata to identify top-$k$ feed keywords

Marilena Oita   Pierre Senellart

- Model a Web page as its DOM (Document Object Level) tree
- Distinguish significant content nodes
  - textual DOM leaf nodes
  - at least one keyword
- Keywords automatically acquired:
  - *Tf-Idf* analysis
    - IN: set of sample Web pages
    - OUT: top-$k$ tf-idf weighted terms
  - basic text preprocessing on feed item metadata to identify top-$k$ feed keywords

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

- Model a Web page as its DOM (Document Object Level) tree
- Distinguish significant content nodes
  - textual DOM leaf nodes
  - at least one keyword
- Keywords automatically acquired:
  - *Tf-Idf* analysis
    - IN: set of sample Web pages
    - OUT: top-$k$ tf-idf weighted terms
  - basic text preprocessing on feed item metadata to identify top-$k$ feed keywords

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

- Model a Web page as its DOM (Document Object Level) tree
- Distinguish significant content nodes
  - textual DOM leaf nodes
  - at least one keyword
- Keywords automatically acquired:
  - *Tf-Idf* analysis
    - IN: set of sample Web pages
    - OUT: top-$k$ tf-idf weighted terms
  - basic text preprocessing on feed item metadata to identify top-$k$ feed keywords

Marilena Oita   Pierre Senellart

```xml
−<rss version="2.0">
  −<channel>
      <title>WORLDMag.com</title>
      <link>http://www.worldmag.com</link>
    +<description></description>
      <pubDate>Sun, 14 Aug 2011 03:20:18 GMT</pubDate>
      <language>en</language>
    +<item></item>
    +<item></item>
    −<item>
        <title>Driver wanted</title>
        <link>http://www.worldmag.com/webextra/18476</link>
        <guid isPermaLink="true">http://www.worldmag.com/webextra/18476</guid>
        <pubDate>Thu, 11 Aug 2011 11:34:01 GMT</pubDate>
        <dc:creator>Joel Hannahs</dc:creator>
      −<description>
          The 'Values Bus' rolls through Iowa in search of a leader on key issues
        </description>
      </item>
    +<item></item>
```

TELECOM
ParisTech

Marilena Oita   Pierre Senellart

DOM element identifier

- tagName
- ⟨attributes⟩
- node index of depth-first search traversal: dfs

Structural pattern $\rightarrow sp_i, i \in 1 : n$

- an *XPath expression*
- describes an identifier of a DOM node which is significant
- example:

$$//div[@id='wrapper' \text{ and } @class='article']$$
$$//p[@dfs=24]$$

DOM element identifier
- tagName
- ⟨attributes⟩
- node index of depth-first search traversal: dfs

Structural pattern $\to sp_i, i \in 1 : n$
- an *XPath expression*
- describes an identifier of a DOM node which is significant
- example:

$$//div[@id='wrapper' \text{ and } @class='article']$$
$$//p[@dfs=24]$$

DOM element identifier

- tagName
- ⟨attributes⟩
- node index of depth-first search traversal: dfs

Structural pattern → $sp_i, i \in 1 : n$

- an *XPath expression*
- describes an identifier of a DOM node which is significant
- example:

$$//div[@id='wrapper' \text{ and } @class='article']$$
$$//p[@dfs=24]$$

TELECOM
ParisTech

SAMPLE PAGE 2

non-significant node

ancestor of a significant node

significant node

SAMPLE PAGE 3

- non-significant node
- ancestor of a significant node
- significant node

## Keyword density

- $x=$ nb. keywords
- $y=$ nb. non-significant terms
- $N=$ total number of terms

$\rightarrow \frac{x}{N}$

## Statistical Corrections:

- $N$ can be small for some nodes
  - *Jeffrey's add-half* estimator $f = \frac{x+1/2}{N+1}$
- $N$ potentially large set $\rightarrow$ margin of error
  - confidence interval of 1
  - one standard deviation (margin of error at 70%) $\sqrt{\frac{f(1-f)}{N}}$

Marilena Oita   Pierre Senellart

## Keyword density

- $x=$ nb. keywords
- $y=$ nb. non-significant terms
- $N=$ total number of terms

$\rightarrow \frac{x}{N}$

## Statistical Corrections:

- $N$ can be small for some nodes
  - *Jeffrey's add-half* estimator $f = \frac{x+1/2}{N+1}$
- $N$ potentially large set $\rightarrow$ margin of error
  - confidence interval of 1
  - one standard deviation (margin of error at 70%) $\sqrt{\frac{f(1-f)}{N}}$

## Keyword density

- $x=$ nb. keywords
- $y=$ nb. non-significant terms
- $N=$ total number of terms

$\rightarrow \frac{x}{N}$

## Statistical Corrections:

- $N$ can be small for some nodes
  - *Jeffrey's add-half* estimator $f = \frac{x+1/2}{N+1}$
- $N$ potentially large set $\rightarrow$ margin of error
  - confidence interval of 1
  - one standard deviation (margin of error at 70%) $\sqrt{\frac{f(1-f)}{N}}$

## Keyword density

- $x =$ nb. keywords
- $y =$ nb. non-significant terms
- $N =$ total number of terms

$\rightarrow \frac{x}{N}$

## Statistical Corrections:

- $N$ can be small for some nodes
  - *Jeffrey's add-half* estimator $f = \frac{x + 1/2}{N + 1}$
- $N$ potentially large set $\rightarrow$ margin of error
  - confidence interval of 1
  - one standard deviation (margin of error at 70%) $\sqrt{\frac{f(1-f)}{N}}$

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

$$f \pm \sqrt{\frac{f(1-f)}{N}} = \frac{x+1/2}{N+1} \pm \frac{1}{N+1}\sqrt{\frac{(x+1/2)\times(y+1/2)}{N}}$$

$$J = \max\left(0, \frac{1}{N+1}\left(x+1/2 - \sqrt{\frac{(x+1/2)\times(y+1/2)}{N}}\right)\right)$$

TELECOM
ParisTech

@DOM node level

- $x=$ keywords
- $y=$ non-significant terms

@Web page level $\rightarrow X, Y$
global context

unexpected content: simpler to describe than to generate

generation complexity $C_w = (x + y) \log (X + Y)$
description complexity $C = x \log X + y \log Y$

$$U = C_w - C$$

Marilena Oita   Pierre Senellart

@DOM node level

- $x=$ keywords
- $y=$ non-significant terms

@Web page level $\rightarrow X$, $Y$
global context

unexpected content: simpler to describe than to generate

generation complexity $C_w = (x + y) \log (X + Y)$
description complexity $C = x \log X + y \log Y$

$$U = C_w - C$$

# Unexpectedness

@DOM node level

- $x=$ keywords
- $y=$ non-significant terms

@Web page level $\rightarrow X, Y$
global context

unexpected content: simpler to describe than to generate

*generation* complexity $C_w = (x + y) \log (X + Y)$
*description* complexity $C = x \log X + y \log Y$

$$U = C_w - C$$

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

# Node Content Informativeness Metric

- measures the interest of a structural pattern $sp_i$
- in a document $d_k$

$$I(sp_i, d_k) = J(sp_i, d_k) \times U(sp_i, d_k)$$

I: informativeness

J: statistical semantic density

U: unexpectedness

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

$$R(sp_i) = \sum_{k=0}^{m} I(sp_i, d_k) \times \text{level}(sp_i) \times \text{nbOcc}_i$$

- Allows ranking generic XPath expressions
  - `//div[@class='wrapper' and (@dfs='27' or @dfs='31')]`

- Final output: subtrees extracted from DOM trees of Web pages

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

$$R(sp_i) = \sum_{k=0}^{m} I(sp_i, d_k) \times \mathrm{level}(sp_i) \times \mathrm{nbOcc}_i$$

- Allows ranking generic XPath expressions
  - `//div[@class='wrapper' and (@dfs='27' or @dfs='31')]`

- Final output: subtrees extracted from DOM trees of Web pages

TELECOM
ParisTech

Marilena Oita    Pierre Senellart

- $\text{FOREST}_{info}$: keywords acquired through tf-idf analysis
- $\text{FOREST}_{feed}$: keywords acquired through feed meta-information

Marilena Oita  Pierre Senellart

# 3 Baselines
# with Alternative Design Choices

- FOREST$_{\text{Cov}}$
  - Same framework as FOREST
  - But score of a pattern is just (normalized) tf-idf weighting of the corresponding DOM nodes

- ABSELEMS
  - Same framework as FOREST$_{\text{info}}$
  - Consider only patterns returning significant leaves, not those that are ancestors of significant leaves

- ABSPATHS
  - Same as ABSELEMS
  - Elements are identified in a pattern by their root-to-leaf path

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

- CETR (WWW, 2010)
  - clustering technique based on a tag ratio per line of the HTML file
  - relies on the fact that text is denser in the main content
- BOILERPIPE (WSDM, 2010)
  - machine learning to determine rules to classify text as content/not-content
  - relies on shallow text features
- DESCRIPTION
  - Main test is just content of the description metadata within a Web feed

- CYAN dataset: dataset used for evaluation of CETR (only 9 different Web sites)
- RED dataset
  - publicly available
  - http://dbweb.enst.fr/software/
  - feed-based (search4Rss); crawl of Web pages referred through feed items
  - manual annotation of the corpus
    - 90 Web sites and 1006 Web pages
    - gold standard: fulltext + metadata (title, author, categories etc.)

TELECOM
ParisTech

- CYAN dataset: dataset used for evaluation of CETR (only 9 different Web sites)
- RED dataset
  - publicly available
  - http://dbweb.enst.fr/software/
  - feed-based (search4Rss); crawl of Web pages referred through feed items
  - manual annotation of the corpus
    - 90 Web sites and 1006 Web pages
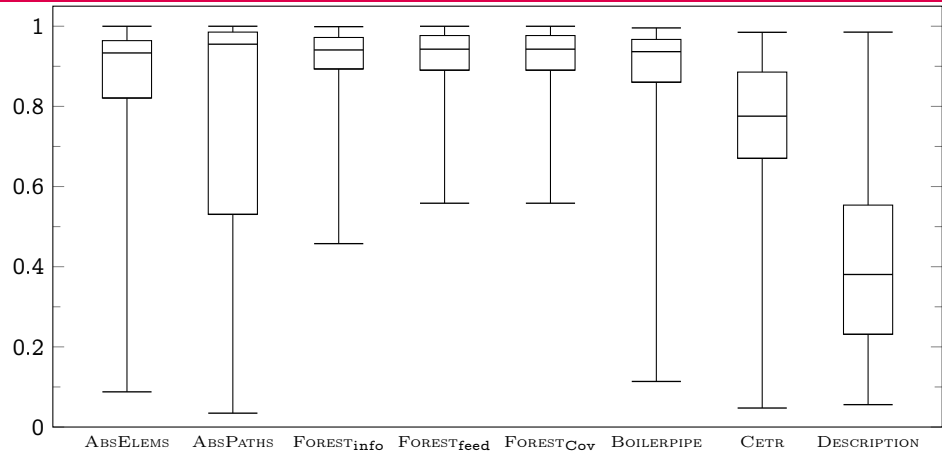    - gold standard: fulltext + metadata (title, author, categories etc.)

TELECOM
ParisTech

| | CYAN | | RED | |
| --- | --- | --- | --- | --- |
| | Prec. | Rec. | Prec. | Rec. |
| ABSELEMS | 65 | 68 | 87 | 93 |
| ABSPATHS | 58 | 64 | 72 | 74 |
| FOREST$_{info}$ | 87 | 99 | 88 | 98 |
| FOREST$_{feed}$ | | | 86 | 98 |
| FOREST$_{Cov}$ | 78 | 89 | 88 | 98 |
| BOILERPIPE | 94 | 97 | 89 | 90 |
| CETR | 65 | 95 | 67 | 93 |
| DESCRIPTION | | | 92 | 31 |

TELECOM
ParisTech

9th and 91th percentile (whiskers), first and third quartile (box) and median (horizontal rule)

Marilena Oita    Pierre Senellart

TELECOM
ParisTech

Marilena Oita   Pierre Senellart

TELECOM
ParisTech

- A fully automatic, robust, effective algorithm for article extraction from dynamically generated Web pages
- Original use of statistical and information-theory–based relevance measures
- Versatile algorithm: has been applied to deep Web object extraction as well (VLDS 2012)
- Requires a source of keywords, which can be external (feed metadata) or internal (informative words on the page itself)
- Extensive and freely available dataset

Merci !

TELECOM
ParisTech

- A fully automatic, robust, effective algorithm for article extraction from dynamically generated Web pages
- Original use of statistical and information-theory–based relevance measures
- Versatile algorithm: has been applied to deep Web object extraction as well (VLDS 2012)
- Requires a source of keywords, which can be external (feed metadata) or internal (informative words on the page itself)
- Extensive and freely available dataset

Merci !

TELECOM
ParisTech